

FACULDADE DE ENGENHARIA DA UNIVERSIDADE DO PORTO

Análise Automática de Melanoma Utilizando Imagens Dermatoscópicas

Bruno Miguel Ferreira Moreira



Mestrado Integrado em Engenharia Informática e Computação

Orientador: Jaime dos Santos Cardoso

Julho de 2017

© Bruno Miguel Ferreira Moreira, 2017

Análise Automática de Melanoma Utilizando Imagens Dermatoscópicas

Bruno Miguel Ferreira Moreira

Mestrado Integrado em Engenharia Informática e Computação

Aprovado em provas públicas pelo Júri:

Presidente: Ademar Aguiar

Vogal Externo: José Costa Pereira

Orientador: Jaime dos Santos Cardoso

18 de Julho de 2017

Resumo

A melanoma é um tipo de cancro da pele que, apesar de ser o menos frequente, é o mais letal na espécie humana [1]. O diagnóstico de melanoma é realizado por um médico especialista – o dermatologista – que observa a pele de um paciente ou a olho nu ou com recurso a um aparelho denominado dermatoscópio.

Em 2012, a WCRFI (acrónimo para *World Cancer Research Fund International*; em português, Fundo Mundial de Pesquisa contra o Cancro) apresentou estatísticas em relação a todos os tipos de cancro ao nível mundial, mostrando que das 14,1 milhões de pessoas diagnosticadas com cancro, 232 mil pessoas estavam diagnosticadas com melanoma [2].

Este tipo de cancro varia muito de aspeto, podendo apresentar-se como uma lesão pigmentada que vai escurecendo, desenvolvendo contornos irregulares ou cores variadas, ao longo do tempo, ou como um nódulo rosa ou encarnado. Pode ser visível em qualquer parte do corpo, sendo o peito e as pernas as zonas mais frequentes. Dado que a melanoma cresce rapidamente, este consegue-se estender para zonas mais internas do corpo, podendo afetar até certos órgãos – sendo esta a razão para que o seu tratamento seja o mais rápido possível.

Esta dissertação tem como objetivo a identificação automática de melanoma em imagens dermatoscópicas de maneira a que o dermatologista possa realizar o tratamento sobre a zona afetada o mais breve possível. Por detrás do algoritmo de identificação, existe um algoritmo de aprendizagem automática que, com base em imagens dermatoscópicas com vários tipos de cancro da pele, contribui para identificar a melanoma eficaz e eficientemente.

Para a construção do algoritmo inteligente, um modelo de características é desenvolvido, através de métodos de aprendizagem automática (como árvores de decisão, máquinas de suporte vetorial, redes neuronais artificiais e entre outros). Para a extração destas características, é necessário obter a lesão da pele da imagem dermatoscópica, através de métodos de segmentação (como o método de Otsu, K-Means, *Watershed* e entre outros). Estas características, juntamente com características próprias do paciente, como a idade ou o sexo, e a identificação de melanoma permitem treinar o algoritmo de aprendizagem. Posto isto, o algoritmo irá prever que classificação atribuir a um novo exame a querer avaliar, classificando-o como a existência de melanoma ou não.

Este algoritmo foi desenvolvido e testado em dados fornecidos e controlados de um concurso de programação denominado *ISIC 2017: Skin Lesion Analysis Towards Melanoma Detection* [3]. Os resultados obtidos sugerem que é necessário aperfeiçoar o sistema antes da sua aplicação na prática clínica.

Palavras-chave: Melanoma; Cancro da Pele; Processamento de Imagem; Aprendizagem Automática

Abstract

Melanoma is a kind of skin cancer that, despite being less frequent, it is the most lethal to human kind. Its diagnosis is done by a specialized doctor – a dermatologist – who observes the patient skin by its own eyes or using a device called dermatoscope [1].

In 2012, WCRFI (acronym for World Cancer Research Fund International) reported statistics for all types of cancer worldwide, showing that of the 14.1 million people diagnosed with cancer, 232.000 people were diagnosed with melanoma [2].

This type of cancer varies in appearance a lot and may present as a pigmented lesion that gets darker, developing irregular contours or variated colors, over time, or as a pink or red nodule. It can be visible on any part of the body, with the chest and legs being the most frequent zones. As melanoma grows, it can spread to the innermost areas of the body, which can affect even certain organs - therefore treatment should be as fast as possible.

This dissertation goal is to give an automatic identification of melanoma in dermatoscopy images so the dermatologist can perform the treatment on the affected area as soon as possible. Behind the identification algorithm, there is a machine learning algorithm that, based on dermatoscopy images with different types of skin cancer, helps to identify melanoma effectively and efficiently.

On building the algorithm, a features' model (extracted by image processing methods) is developed by machine learning methods (like decision trees, support vector machines, artificial neural networks and others). To extract these features, it is necessary to get the skin lesion from the dermatoscopy image, through segmentation methods (such as the Otsu's method, K-Means, Watershed and others). These characteristics along with patient's own characteristics such as age or sex and the identification of melanoma (or not) allow to train the algorithm. In this way, the algorithm will predict which classification to assign to a new exam to evaluate, classifying it as the existence of melanoma or not.

This algorithm was developed and tested using controlled data given by the contest *ISIC 2017: Skin Lesion Analysis Towards Melanoma Detection* [3]. The results suggest that it's necessary to improve the system before its application on clinical practice.

Key-words: Melanoma; Skin Cancer; Image Processing; Machine Learning

Agradecimentos

Esta dissertação foi realizada entre o dia 18 de Abril de 2017 e o dia 27 de Junho de 2017, o que equivale a cerca de dois meses e meio de trabalho, e corresponde a uma alternativa de dissertação proposta pelo meu orientador quando a minha dissertação inicial, proposta em Setembro, não tinha os recursos necessários para a realizar e a tempo devido.

As áreas em que esta dissertação se enquadrava, nomeadamente a dermatologia, o processamento de imagem e a aprendizagem automática, tornaram-se desafios dada a minha inexperiência nos temas em si. Por esta mesma inexperiência, eu escolhi esta dissertação para aprender, melhorar e crescer.

Como uma velha frase diria, a gratidão é o único tesouro dos humildes. E não deixa de ser verdade.

Dado o constrangimento de tempo e a minha dificuldade nas áreas, gostaria de agradecer principalmente ao professor Jaime Cardoso, pelo apoio e paciência nestes oito meses em que me orientou, especialmente nestes últimos três meses; pois se não fosse o meu orientador, talvez esta dissertação não estaria escrita.

Queria agradecer igualmente a vários amigos que me ajudaram e ensinaram quase que um novo mundo de conceitos, especialmente na área de saúde, quer na dissertação anterior quer nesta dissertação. Mas assim seria também injusto não agradecer aqueles que estiveram comigo ao longo destes cinco anos de faculdade, pois foram com eles com quem fui à luta nas provas e constrangimentos mais difíceis.

E principalmente não podia deixar de agradecer à minha família – aos meus pais, ao meu irmão, aos meus avós, aos meus tios e aos meus primos, pois toda a família teve paciência comigo e pude apoiar-me nos bons e nos maus momentos, mesmo quando tudo parecia perdido, mesmo quando outras pessoas me queriam prejudicar, mesmo quando muita gente me dizia que nunca seria capaz de chegar longe ou de ser alguém.

A toda a família, agradeço-lhes de verdade. Obrigado.

Bruno Miguel Ferreira Moreira

Deleita-te com a tua obra e aprecia a obra dos outros, pois quem não tem obras não tem nada para mostrar.

O meu pai.

Conteúdo

Introdução.....	1
1.1 Contexto	1
1.2 Motivação e objetivos	1
1.3 Resultados alcançados.....	2
1.4 Estrutura da Dissertação.....	3
Estado da Arte.....	5
2.1 Introdução	5
2.2 Pele Humana	6
2.2.1 Anatomia	6
2.2.2 Lesões da Pele	7
2.2.3 Métodos de Identificação	9
2.3 Aprendizagem Automática.....	10
2.3.1 Árvore de Decisão	11
2.3.2 Rede Neuronal Artificial	12
2.3.3 Máquina de Suporte Vetorial	14
2.3.4 Floresta Aleatória	15
2.3.5 Naive Bayes	16
2.3.6 K-Means.....	18
2.3.7 GMM.....	19
2.4 Processamento de Imagem	21
2.4.1 Modelos de Cor	21
2.4.2 Segmentação	21
2.4.2.1 Método de <i>Thresholding</i>	22
2.4.2.2 Método de Otsu	22
2.4.2.3 K-Means	23
2.4.2.4 GMM	23
2.4.2.5 Watershed	23
2.4.2.6 Cortes Normalizados	24
2.4.2.7 Closed Shortest Path	26

2.4.3	Extração de Características	28
2.4.3.1	Propriedades Geométricas	29
2.4.3.2	Assimetria de um objeto	30
2.4.3.3	Propriedades da Cor	31
2.4.3.4	Matriz de Coocorrências	31
2.5	Métricas.....	33
2.6	Trabalhos Relacionados	36
Metodologia		39
3.1	Introdução	39
3.2	Análise do Problema	39
3.3	Proposta de Solução	40
3.4	Proposta de Validação.....	42
Implementação		45
4.1	Introdução	45
4.2	Arquitetura do Projeto.....	45
4.2.1	Arquitetura Física.....	45
4.2.2	Arquitetura Tecnológica.....	45
4.2.3	Estrutura do Projeto.....	47
4.3	Desenvolvimento.....	48
4.3.1	Solução.....	48
4.3.2	Validação.....	50
4.4	Resultados	51
4.4.1	Segmentação	51
4.4.2	Classificação	52
Conclusão.....		55
Referências.....		57

Lista de Figuras

Figura 1: Anatomia da Pele Humana. Fonte: [6]	6
Figura 2: Localização dos melanócitos e da melanina na epiderme. Fonte: [7]	7
Figura 3: Exemplo de ceratose seborreica.	8
Figura 4: Exemplo de nevo.	8
Figura 5: Exemplo de melanoma.	8
Figura 6: Aplicação do dermatoscópio na identificação de lesões da pele. Fonte: [12]	9
Figura 7: Árvore de Decisão para o exemplo dado. Fonte: [19]	12
Figura 8: Arquitetura de uma rede neuronal artificial. Fonte: [22]	13
Figura 9: Exemplo de uma SVM. Fonte: [28]	14
Figura 10: Aplicação de uma função de nuclearização. Fonte: [29]	15
Figura 11: Exemplo de uma floresta aleatória. Fonte: [32]	16
Figura 12: Função de thresholding. Fonte: [49]	22
Figura 13: Preenchimento da imagem raiz pelo algoritmo de Watershed.	24
Figura 14: Exemplo da linearização com base numa vizinhança de raio 5. Fonte: [61]	27
Figura 15: Exemplo da obtenção dos segmentos de reta.	31
Figura 16: Exemplo de duas matrizes de ocorrência geradas a partir de duas imagens.	32
Figura 17: Esquema da Fase de Treino.	41
Figura 18: Esquema da Fase de Validação.	42
Figura 19: Esquema da Fase de Aplicação.	42
Figura 20: Máscara obtida pelo método de Otsu.	43
Figura 21: Máscara de referência.	43
Figura 22: Esquema do programa de segmentação.	48
Figura 23: Esquema do programa de extração de características.	49
Figura 24: Esquema do programa de construção do modelo.	49
Figura 25: Esquema do programa de previsão de classificações.	50

Lista de Tabelas

Tabela 1: Exemplo de uma tabela de frequência para o atributo Estado do Tempo.	
Fonte: [36]	17
Tabela 2: Exemplo de uma tabela de probabilidades para o atributo Estado do Tempo.	
Fonte: [36]	18
Tabela 3: Categorias de características relacionadas com objetos segmentados.	28
Tabela 4: Matriz de Confusão.	33
Tabela 5: Versões dos módulos de Python instalados e utilizados no projeto.	46
Tabela 6: Métricas obtidas pelos métodos de segmentação.	51
Tabela 7: Matriz de Confusão para um modelo de cada base de segmentação.	52
Tabela 8: Métricas calculadas para cada modelo de uma base de segmentação.	53

Lista de Equações

Equação 1: Equação da Entropia no contexto das Árvores de Decisão.	11
Equação 2: Equação da Entropia para um Atributo no contexto das Árvores de Decisão.	12
Equação 3: Equação do Ganho no contexto das Árvores de Decisão.	12
Equação 4: Fórmula da função logística.	13
Equação 5: Teorema de Bayes.	16
Equação 6: Teorema de Bayes utilizando a terminologia da probabilidade bayesiana.	17
Equação 7: Teorema de Bayes aplicado ao problema de classificação.	17
Equação 8: Teorema de Bayes e a probabilidade em cadeia.	17
Equação 9: Fórmula para obter observações mais próximas de cada centroide.	18
Equação 10: Fórmula para modificar cada centroide.	19
Equação 11: Fórmula da probabilidade de um GMM.	19
Equação 12: Fórmula para a função de distribuição gaussiana.	19
Equação 13: Fórmula da probabilidade de um GMM para múltiplos atributos.	19
Equação 14: Fórmula para a função de distribuição gaussiana para múltiplos atributos.	20
Equação 15: Fórmula para inicializar as variâncias.	20
Equação 16: Fórmula da probabilidade de um elemento i pertencer a uma componente k .	20
Equação 17: Fórmula para atualizar os pesos da fórmula do algoritmo GMM.	20
Equação 18: Fórmula para atualizar as médias da fórmula do algoritmo GMM.	20
Equação 19: Fórmula para atualizar os desvios-padrões da fórmula do algoritmo GMM.	21
Equação 20: Fórmula da obtenção do threshold. Fonte: [51]	23
Equação 21: Fórmula para o corte em A e B.	24
Equação 22: Fórmula para definir os pesos das ligações do grafo. Fonte: [58]	24
Equação 23: Fórmula para calcular a medida de dissociação.	25
Equação 24: Fórmula para calcular a medida de associação.	25
Equação 25: Relação entre a medida de associação e a medida de dissociação.	25
Equação 26: Definição da matriz diagonal.	25

Equação 27: Nova fórmula do corte para o algoritmo do corte normalizado.	25
Equação 28: Equação para minimizar o corte normalizado, calculando os vetores próprios.	26
Equação 29: Definição do peso de uma aresta que parte de um vértice.	27
Equação 30: Fórmula do diâmetro equivalente.	29
Equação 31: Fórmula da compacidade.	29
Equação 32: Fórmula da circularidade.	29
Equação 33: Fórmula da solidez.	29
Equação 34: Fórmula da retangularidade.	30
Equação 35: Fórmula da proporção.	30
Equação 36: Fórmula da excentricidade.	30
Equação 37: Fórmula do momento ij .	30
Equação 38: Fórmula da Inclinação.	31
Equação 39: Fórmula do segundo momento angular.	32
Equação 40: Fórmula do contraste.	32
Equação 41: Fórmula da correlação.	32
Equação 42: Equação da Sensibilidade.	34
Equação 43: Equação da Especificidade.	34
Equação 44: Equação da Queda.	34
Equação 45: Equação da Taxa de Falha.	34
Equação 46: Equação da Precisão.	35
Equação 47: Equação do Valor de Previsão Negativa.	35
Equação 48: Equação da Taxa de Descobrimiento Falso.	35
Equação 49: Equação da Taxa de Omissão Falsa.	35
Equação 50: Equação da Exatidão.	35
Equação 51: Equação da Prevalência.	35
Equação 52: Equação do Índice de Jaccard.	36
Equação 53: Equação do Índice de Sørensen-Dice.	36

Abreviaturas e Símbolos

WCRFI	<i>World Cancer Research Fund International</i>
IEEE	<i>Institute of Electrical and Electronics Engineers</i>
ISBI	<i>International Symposium on Biomedical Imaging</i>
ISIC	<i>International Skin Imaging Collaboration</i>
CSV	<i>Comma-separated values</i>
DT	<i>Decision Tree</i>
ANN	<i>Artificial Neural Network</i>
SVM	<i>Support Vector Machine</i>
RF	<i>Random Forest</i>
NB	<i>Naïve Bayes</i>
GMM	<i>Gaussian Mixture Model</i>
EM	<i>Expectation-Maximization</i>
DAG	<i>Directed acyclic graph</i>
DFS	<i>Deep-First Search</i>
RGB	<i>Red Green Blue</i>
BGR	<i>Blue Green Red</i>
HSV	<i>Hue Saturation Value</i>
NCut	<i>Normalized Cut</i>
NP	<i>Non-Deterministic</i>
VP	Verdadeiro Positivo
FP	Falso Positivo
FN	Falso Negativo
VN	Verdadeiro Negativo
JSON	<i>JavaScript Object Notation</i>
RAM	<i>Random Access Memory</i>
PKL	<i>PacKing List</i>

Capítulo 1

Introdução

1.1 Contexto

Numa sociedade onde a tecnologia tem um papel tão importante na vida das pessoas, a sua aplicação na saúde e no bem-estar deve ser uma prioridade fundamental. Com o crescimento da população mundial, há cada vez mais pessoas a recorrer a institutos de saúde, nomeadamente a hospitais, centros de saúde ou clínicas de saúde para receber tratamento. Este tratamento é efetuado por médicos que cada vez mais recorrem à tecnologia com o objetivo de serem auxiliados na sua tarefa.

O tempo é uma dádiva, a vida é uma dádiva. É essencial preservar a vida de uma pessoa, custo o que custar. A pele é o órgão mais extenso e exterior do ser humano; é o nosso retrato, aquele que permite outras pessoas nos identificarem. Quando a pele fica afetada por ferimentos ou danos, por doenças ou por anomalias, é vital que um médico especialista – o dermatologista – atue de imediato para salvaguardar o paciente.

Um dermatoscópio é um dos instrumentos não invasivos que permite identificar anomalias e/ou lesões na pele. As imagens que se obtêm através do uso deste mesmo aparelho – as imagens dermatoscópicas – são utilizadas para que o médico possa identificar qualquer anomalia com melhor precisão, até mesmo ter um histórico da zona da pele em análise.

1.2 Motivação e objetivos

Quando se trata de identificação de anomalias, há que ter em conta que certas anomalias podem ser fatais para o ser humano, tais como a melanoma. A necessidade de intervir o mais breve possível torna-se uma prioridade, não uma opção – é importante que o tempo que se demora a identificar uma anomalia seja o menor possível.

O objetivo desta dissertação é utilizar algoritmos de aprendizagem automática capazes de gerar modelos de características obtidas através de métodos de processamento de imagem sobre imagens dermatoscópicas, com o intuito de não só identificar melanoma em novos casos em tempo reduzido como também auxiliar os dermatologistas que tenham maior dificuldade de análise, especialmente os menos experientes.

Para o desenvolvimento do algoritmo de identificação de melanoma, recorreu-se a um desafio online criado e patrocinado por um ramo da **IEEE** (acrónimo para *Institute of Electrical and Electronics Engineers*; em português, Instituto de Engenheiros Eletricistas e Eletrónicos), de maneira a utilizar dados fiáveis e controlados. A **IEEE International Symposium on Biomedical Imaging** (ISBI como acrónimo) é uma conferência científica dedicada aos aspetos matemáticos, algorítmicos e computacionais da imagem biomédica, com o intuito de partilhar experiência com as várias comunidades existentes e contribui para uma abordagem integradora da imagem biomédica [4].

Em 2017, a ISBI criou um desafio denominado **Análise de Lesão da Pele para a Detecção de Melanoma** (em inglês, *Skin Lesion Analysis Towards Melanoma Detection*) em parceria com a **International Skin Imaging Collaboration** (ISIC como acrónimo; em português, Colaboração Internacional de Imagem de Pele), cuja intenção desta colaboração é melhorar o diagnóstico da melanoma [3][5].

O objetivo do desafio é aplicar algoritmos de processamento de imagem e aprendizagem automática para detetar **melanoma**, **nevo** e **ceratose seborreica** num conjunto de imagens dermatoscópicas.

Este desafio é composto por três fases: **segmentação da lesão**, **extração de características da lesão** e **classificação da lesão**. Para tal, são fornecidas 2000 imagens dermatoscópicas como dados de modelação (tanto de treino como de validação). Juntamente com as imagens é fornecido um ficheiro CSV (acrónimo para *Comma-separated values*; em português, Valores separados por vírgula) que contem a idade e o sexo do paciente ao qual foi obtida a imagem dermatoscópica.

O ISIC 2017 teve início no dia 9 de Dezembro de 2016 e fim no dia 1 de Março de 2017.

1.3 Resultados alcançados

Foi realizado um estudo exaustivo dos métodos de segmentação propostos, aplicando todos os métodos a todas as imagens dermatoscópicas da base de dados. Das métricas calculadas das imagens segmentadas para com as imagens segmentadas de referência ou do especialista, conclui-se que o algoritmo **Closed Shortest Path** é aquele que tem melhor sensibilidade e exatidão. Sendo assim, definiu-se este algoritmo como aquele que realiza a segmentação da imagem.

De seguida, foram extraídas as características propostas, com base num número selecionado de imagens segmentadas pelo algoritmo *Closed Shortest Path* e de imagens segmentadas de referência e de seguida guardadas em ficheiros. Vários algoritmos de aprendizagem automática foram aplicados as características extraídas, gerando os modelos de características respetivos.

Uma vez gerados os modelos, estes mesmos foram utilizados para prever a classificação a outro conjunto de imagens da base de dados. As métricas calculadas permitiram concluir que as características extraídas ou não são suficientes ou a sua qualidade é reduzida para classificar, dado que as redes neuronais e as máquinas de suporte vetorial não conseguiram identificar melanoma em nenhum caso do conjunto de validação, apesar de ambos os algoritmos de aprendizagem automática terem a maior exatidão, mas inferior a 75%. Para o modelo gerado a partir do algoritmo de segmentação, as **Florestas Aleatórias** são o modelo a aplicar no algoritmo final, dado que é o algoritmo com maior exatidão (67%).

1.4 Estrutura da Dissertação

Este documento está estruturado em cinco capítulos. Uma vez concluída esta introdução do tema da dissertação, o Capítulo 2 focará no estado da arte. Posteriormente, o Capítulo 3 será abordada a especificação do projeto, onde será discutida a solução a atacar para o problema. No Capítulo 4 será abordada a implementação da solução mencionada no capítulo anterior e exposto os resultados obtidos. Por fim, no Capítulo 5 serão realçadas as conclusões e o possível trabalho futuro.

Capítulo 2

Estado da Arte

2.1 Introdução

Neste capítulo, será abordado o estado da arte relacionada com a dissertação; isto é, será levantado um estudo sobre os fundamentos da dissertação. Primeiro, será abordada a pele humana, nomeadamente a sua anatomia, o cancro da pele e os vários métodos de identificação do cancro da pele. Posteriormente serão abordados tanto os vários algoritmos de aprendizagem automática como as várias técnicas de processamento de imagem, nomeadamente os vários métodos de segmentação e extração de características.

De seguida, será analisado as métricas utilizadas para medir o desempenho dos vários métodos referidos anteriormente.

Por fim, serão abordados os diversos trabalhos já realizados no contexto da dissertação e que contribuíram para o desenvolvimento de diversas soluções para o mesmo tema.

2.2 Pele Humana

2.2.1 Anatomia

A pele humana (em inglês, *human skin*) é o órgão mais largo e o que cobre a maioria do corpo humano. A pele, existente nos mamíferos, ajuda a proteger o organismo de micróbios, ajuda a regular a temperatura corporal e permite sensações como tato, calor e frio [6].

A pele é composta essencialmente por três camadas:

- a epiderme (em inglês, *epidermis*), a camada mais externa da pele; funciona como uma barreira impermeável e protetora de luz solar. Oferece o tom que define a pele, este graças a células denominadas melanócitos (em inglês, *melanocytes*) que produzem pigmentos de melanina (em inglês, *melanin*).
- a derme (em inglês, *dermis*), a camada mais interna que a epiderme; contem tecidos (em inglês, *tissues*) fortemente ligados, foliculos capilares (em inglês, *hair follicles*) e glândulas sudoríparas (em inglês, *sweat glands*).
- a hipoderme (em inglês, *hypodermis*) ou tecido celular subcutâneo (em inglês, *subcutaneous tissue*), a camada mais interna da pele; é constituída essencialmente por tecidos, gordura (em inglês, *fat*) e por vasos sanguíneos (em inglês, *blood vessels*).

A Figura 1 ilustra a anatomia da pele humana, da mesma maneira que na Figura 2 é possível localizar os melanócitos e a melanina na epiderme.

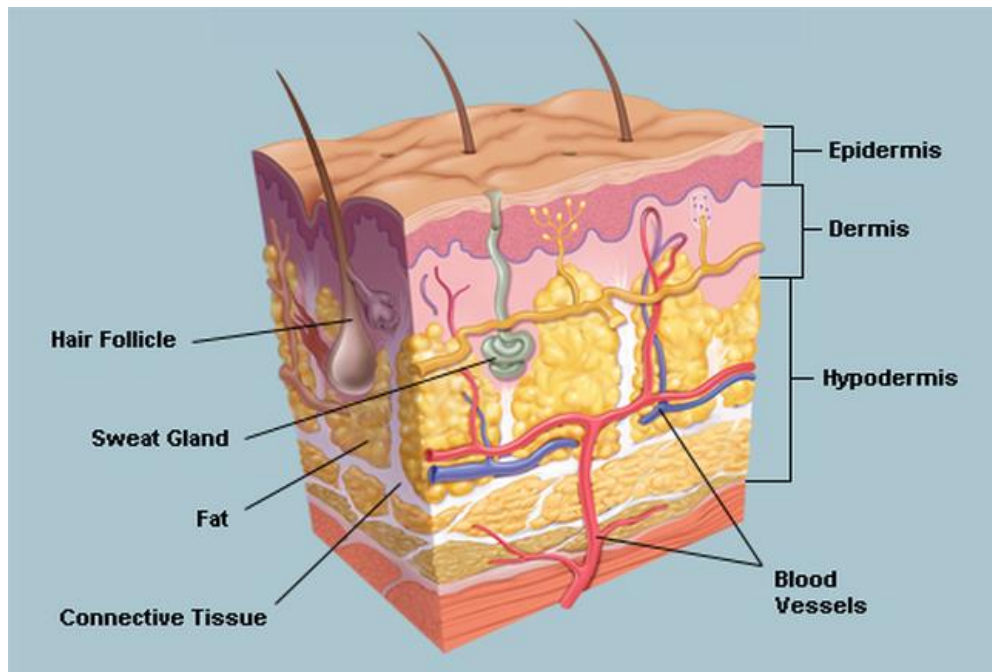


Figura 1: Anatomia da Pele Humana. Fonte: [6]

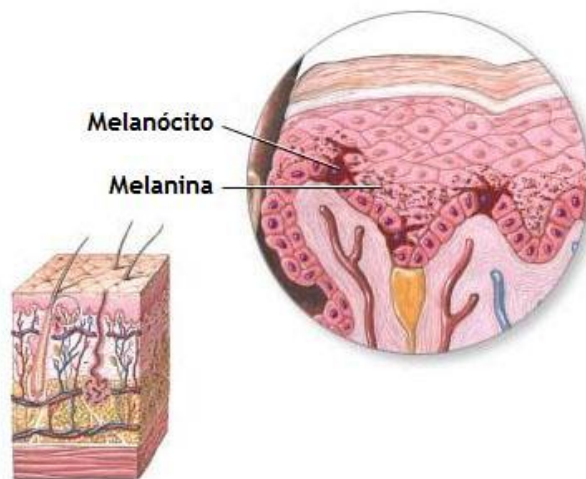


Figura 2: Localização dos melanócitos e da melanina na epiderme. Fonte: [7]

2.2.2 Lesões da Pele

Neoplasia (em inglês, *neoplasm*) é um crescimento anormal de um tecido. Quando por vezes o crescimento gera uma massa, essa massa é designada de tumor (em inglês americano, *tumor*; em inglês britânico, *tumour*).

Existem quatro grupos principais de neoplasia, definidas pela Organização Mundial de Saúde [8]:

- Neoplasia benignas;
- Neoplasia potencialmente maligna;
- Neoplasia maligna – também conhecida como câncer (em inglês, *cancer*);
- Neoplasia de comportamento incerto ou desconhecido.

Uma lesão da pele (em inglês, *skin lesion*) é um termo que se refere a qualquer anormalidade na pele. Estas lesões têm origem tanto por infecções, alergias ou absorção abusiva de raios ultravioleta como por manifestações de neoplasias. Existe uma diversidade de lesões de pele [9], mas o foco desta dissertação está nas seguintes:

- ceratose seborreica ou queratose seborreica (em inglês, *seborrheic keratosis*), que é uma lesão da pele que é considerada como um tumor benigno; tem uma aparência que se assemelha a uma verruga [10][11].
- nevo (em inglês, *nevus*), que é uma lesão da pele que é considerada como um tumor benigno;
- melanoma, que é uma lesão da pele que é considerada como um tumor maligno (ou seja, câncer); tanto pode surgir como uma evolução grave de um nevo como pode surgir sem haver um nevo previamente. Apresenta ser uma lesão pigmentada que vai escurecendo, desenvolvendo contornos irregulares ou cores variadas, ao longo do tempo, ou como um nódulo rosa ou encarnado.

As lesões da pele detalhadas anteriormente têm como exemplo, respectivamente, as Figura 3, Figura 4 e Figura 5. As lesões identificadas como tumores benignos podem ser facilmente removidas da pele, enquanto que a melanoma é um caso mais grave, em que quanto mais tempo se perde, mais fatal pode ser para a pessoa em causa.

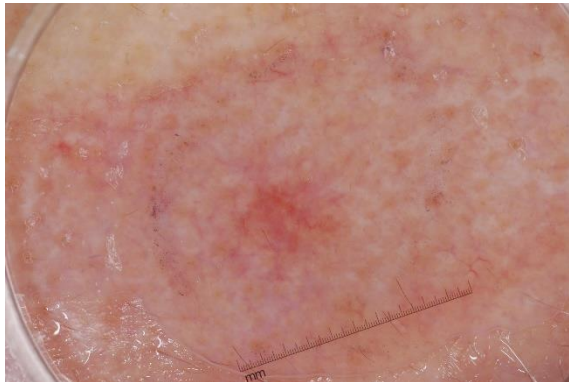


Figura 3: Exemplo de ceratose seborreica.

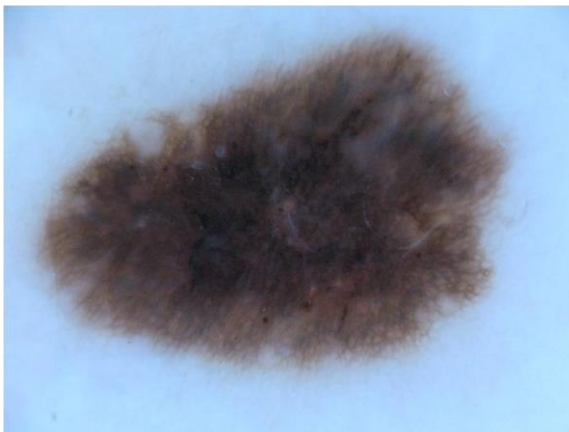


Figura 4: Exemplo de nevo.

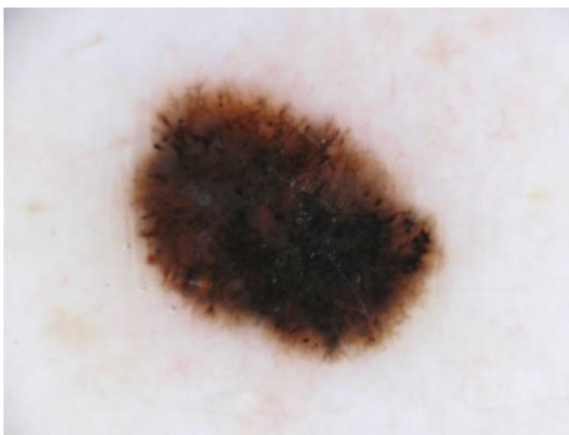


Figura 5: Exemplo de melanoma.

2.2.3 Métodos de Identificação

Existem uma diversidade de métodos de identificação de lesões da pele. Na categoria da capacidade de visualização, existe duas categorias principais:

- tradicional, em que o dermatologista identifica a olho nu as lesões da pele;
- tecnológico, em que o dermatologista recorre a aparelhos que o auxiliam.

Um dos aparelhos mais utilizados são os dermatoscópios (em inglês, *dermatoscopes*), como é o exemplo da Figura 6. Com as imagens obtidas destes aparelhos – denominadas imagens dermatoscópicas (em inglês, *dermoscopic images*) – é possível utilizar regras para a identificação de possíveis lesões.



Figura 6: Aplicação do dermatoscópio na identificação de lesões da pele. Fonte: [12]

Como a melanoma é o foco da dissertação, apenas se vai restringir a regras de identificação de melanoma. A principal regra e a mais conhecida é a regra **ABCDE** [13][14]:

- assimetria (em inglês, *asymmetry*): a lesão é assimétrica?
- bordo (em inglês, *border*): o bordo da lesão é irregular?
- cor (em inglês americano, *color*; em inglês britânico, *colour*): existem pigmentos, picos e vales de cor
- diâmetro (em inglês, *diameter*): o diâmetro é superior a 6mm?
- evolução (em inglês, *evolution*): o tamanho, a forma e o comportamento alterou com o tempo?

Importante afirmar que às vezes a regra ABCDE pode ser abreviada para ABC, pois a maioria dos casos a identificação de melanoma é a primeira vez numa dada região do corpo. Por isso mesmo, não faz muito sentido pensar na evolução da lesão para os novos pacientes. Quando as imagens não dão informação sobre a escala da mesma, obtida pelo dermatoscópio, o diâmetro pode ser descartado.

2.3 Aprendizagem Automática

Aprendizagem automática (em inglês, *machine learning*) é, segundo Arthur Samuel, definido como “a área de estudo que permite dar aos computadores a capacidade de aprender sem serem programados explicitamente” [15].

Apesar de ser uma definição um pouco informal, Tom Mitchell elaborou uma definição mais moderna: "Um programa aprende da experiência E em relação a um conjunto de tarefas T e medidas de desempenho P, se o seu desempenho nas tarefas de T, desempenho medido por P, melhora com a experiência E" [16].

Os algoritmos de aprendizagem automática podem agrupados em três categorias, segundo o tipo de tarefas a executar:

- **Aprendizagem supervisionada** (em inglês, *supervised learning*), em que se fornece um conjunto de N elementos de entrada ao algoritmo e sabe-se exatamente qual o tipo de saída para cada um dos N elementos.
- **Aprendizagem não-supervisionada** (em inglês, *unsupervised learning*), em que se fornece um conjunto de N elementos de entrada ao algoritmo, mas não se sabe o tipo de saída para cada um dos N elementos.
- **Aprendizagem por reforço** (em inglês, *reinforcement learning*), em que se fornece um conjunto de estados E sobre o meio-ambiente que o computador está inserido, um conjunto de ações A em que o computador pode executar de modo a transitar entre os vários estados e um conjunto de recompensas/penalizações R que se obtém ao tomar uma ação de A. O algoritmo tem que ser capaz de prever as melhores ações para o meio ambiente que se encontra.

Para cada categoria acima, os algoritmos de aprendizagem automática podem ser agrupados em três categorias, segundo o tipo de dados de saída:

- **Classificação** (em inglês, *classification*), em que o algoritmo tem que prever um valor discreto.
- **Regressão** (em inglês, *regression*), em que o algoritmo tem que prever um valor contínuo.
- **Agrupamento** (em inglês, *clustering*), em que o algoritmo tem que separar um conjunto N de dados de entrada em M conjuntos de dados de saída não conhecidos, em que $M < N$. Faz parte unicamente dos algoritmos de aprendizagem não-supervisionada.

A aprendizagem supervisionada é constituída por duas fases: a **fase de modelação** (em inglês, *development phase*) e a **fase de aplicação** (em inglês, *application phase*). A fase de modelação, por si só, pode ser decomposta em mais duas fases: a **fase de treino** (em inglês, *training phase*) e a **fase de validação** (em inglês, *validation phase*).

Na fase de modelação, é fornecido um conjunto de dados M composto por X observações e Y resultados das respetivas observações. Este conjunto de dados é dividido em duas partes, uma maioria T para a fase de treino e a outra minoria V para a fase de validação – tipicamente, um terço ou um quarto dos dados são utilizados na fase de validação.

Na fase de treino, o objetivo é construir o modelo de dados com base no conjunto de dados T e com base no algoritmo de aprendizagem em causa. Na fase de validação, o objetivo é testar/validar o modelo de dados construído. Para tal, o algoritmo vai tentar prever um conjunto de resultados Y_V' para as observações X_V do conjunto de dados V , comparando com os resultados reais Y_V das respetivas observações.

Na fase de aplicação, o objetivo é utilizar o modelo construído nos desafios com que se enquadra o mesmo, prevendo e afirmando resultados para novas observações.

Nas seguintes subseções, serão abordados vários algoritmos de aprendizagem automática.

2.3.1 Árvore de Decisão

Uma **árvore de decisão** (em inglês, *decision tree*; DT como acrónimo) é um algoritmo de aprendizagem supervisionada, que pode ser considerada tanto de classificação como de regressão. O objetivo desta estrutura é expressar regras capazes de fornecer uma decisão a tomar com base numa observação, regras essas representadas e definidas sob a forma de uma árvore. Vários algoritmos que implementam árvores de decisão já foram desenvolvidos, nomeadamente ID3, C4.5 e CART [17][18][19][20][21].

Cada nó da árvore representa uma regra, uma condição (sobre uma observação) a verificar. Cada ramo de um nodo representa o resultado esperado da verificação da regra. E por fim, cada folha de um ramo representa a decisão final. Para a construção de uma árvore de decisão, vários conceitos são aplicados, nomeadamente a **entropia** (em inglês, *entropy*) e a **informação de ganho** (em inglês, *information gain*) ou simplesmente **ganho**. A entropia encontra-se traduzida nas Equação 1 e Equação 2, enquanto que o ganho se encontra traduzido na Equação 3.

$$H(S) = - \sum_{i=1}^N [p_i * \log_2(p_i)]$$

Equação 1: Equação da Entropia no contexto das Árvores de Decisão.

Na Equação 1, S é o conjunto de dados composto por X atributos independentes (as observações) e pelo atributo dependente Y (a decisão); Y é composto por N classes. Na construção da raiz da árvore considera-se S o conjunto total de dados, e por isso p_i é a probabilidade de ocorrência de uma classe de Y em S . Quando se considera um ramo, dados de S relativos à regra criada no nodo do ramo são ignorados para a nova iteração, gerando um novo conjunto S' . Sendo assim, a equação $H(S')$ pode ser igualmente escrita por $H(\text{Atributo} = \text{valor})$, em que p_i é a probabilidade de ocorrência de um atributo ser igual ao valor em S' .

$$H_{X_n}(S) = \sum_{i=1}^M \left[\frac{|S_i|}{|S|} * H(S_i) \right]$$

Equação 2: Equação da Entropia para um Atributo no contexto das Árvores de Decisão.

Na Equação 2, X_n é uma observação a calcular a sua entropia em S , composto por M valores e S_i é o conjunto de dados em S que é composto por observações em que X_n é igual a um valor i .

$$G_{X_n}(S) = H(S) - H_X(S)$$

Equação 3: Equação do Ganho no contexto das Árvores de Decisão.

Para um dado nodo, uma vez calculada os ganhos para cada atributo X_n , a regra para o nodo é aquela cujo $G_{X_n}(S)$ é o maior. A Figura 7 ilustra um exemplo de uma árvore de decisão.

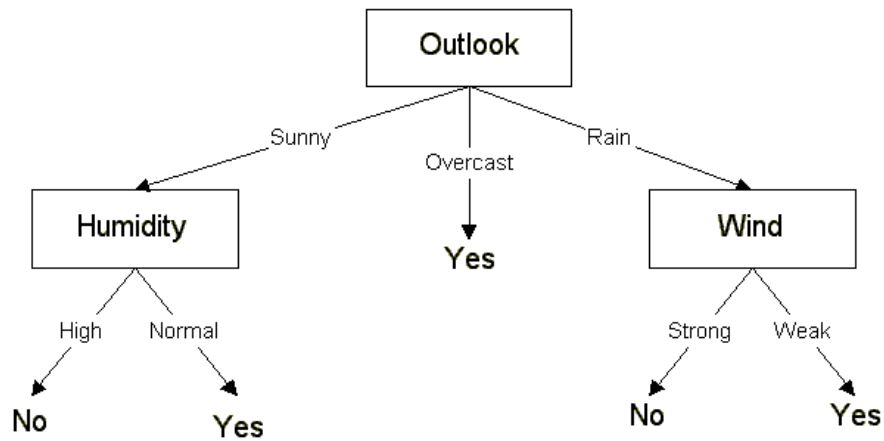


Figura 7: Árvore de Decisão para o exemplo dado. Fonte: [19]

2.3.2 Rede Neuronal Artificial

Uma **rede neuronal artificial** (em inglês, *artificial neural network*; ANN como acrónimo) é um algoritmo tanto de aprendizagem supervisionada, que pode ser considerada tanto de classificação como de regressão, como de aprendizagem não-supervisionada [23].

A rede em si é um grafo. Na Figura 8 encontra-se um exemplo da arquitetura de uma rede neuronal, onde se pode visualizar a existência de três camadas principais:

- a **camada de entrada** (em inglês, *input layer*), em que se fornece os dados de entrada ou as observações X .
- as **camadas ocultas** ou intermédias (em inglês, *hidden layers*)

- a **camada de saída** (em inglês, *output layer*), em que se fornece os dados de saída ou classificações Y, caso seja um problema de aprendizagem supervisionada; independentemente, é nesta camada onde se obtêm as previsões.

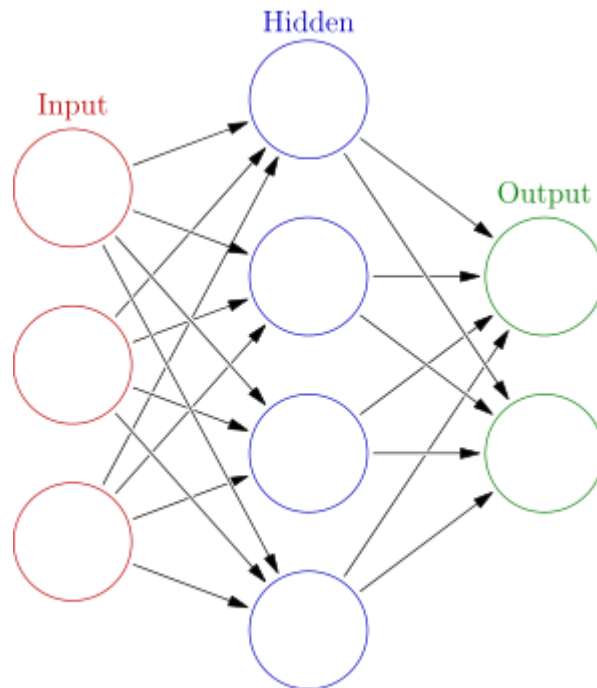


Figura 8: Arquitetura de uma rede neuronal artificial. Fonte: [22]

As camadas encontram-se interligadas por uma conexão entre dois neurónios e cada conexão tem um respetivo peso P. Cada neurónio é constituído por dois valores: valor de receção R e o valor de ativação A. O valor de ativação de um neurónio é obtido através da aplicação de uma função de ativação f sobre o valor de receção, em que tipicamente a função de ativação é a função de logística (representada na Equação 4). O valor de receção é dado como a média ponderada entre os valores de ativação dos neurónios antecedentes e os respetivos pesos das conexões.

$$f(x) = \frac{1}{1 + e^{-x}}$$

Equação 4: Fórmula da função logística.

Quando numa rede neuronal existem ligações entre os neurónios capazes de gerar ciclos na rede, essa rede é denominada de **rede neuronal recorrente** (em inglês, *recurrent neural network*). Quando não se verifica nenhum ciclo numa rede neuronal, essa rede é denominada de **rede neuronal de alimentação progressiva** (em inglês, *feed-forward neural network*). No caso particular em que se pretende classificar uma observação de um dado tipo, a rede é denominada de **perceptron** [24].

A construção de um perceptron é simples: realizar várias iterações de realimentação e retropropagação (em inglês, *backpropagation*), nesta precisa ordem. Na realimentação, o objetivo é propagar os dados de entrada em cada neurônio até aos neurônios de saída. Na retropropagação (no caso particular da aprendizagem supervisionada), o objetivo é corrigir os pesos das conexões de modo a que os valores de ativação dos neurônios da camada de saída se aproximem dos dados de saída ou classificações fornecidas. Na previsão de resultados, uma simples realimentação é suficiente para obter os resultados previstos nos neurônios da camada de saída.

2.3.3 Máquina de Suporte Vetorial

Uma **máquina de suporte vetorial** (em inglês, *support vector machine*; SVM como acrónimo) é um algoritmo de aprendizagem supervisionada, que pode ser considerada tanto de classificação como de regressão.

Dado um conjunto de observações X , compostas por N atributos, e as respetivas classificações Y , o objetivo é saber qual o **hiperplano**, de dimensão N , que melhor separa as observações X em relação a diferentes classificações Y [25][26].

No caso particular de uma classificação binária (ou seja, duas classes) em que as observações têm apenas 2 atributos (exemplo esse ilustrado na Figura 9), o hiperplano encontra-se representado segundo uma linha: a obtenção deste hiperplano é realizada a partir do cálculo de uma distância mínima – denominada **margem** (em inglês, *margin*) – entre um ou mais elementos das respetivas duas classes. A linha separa a região dentro da margem a meio. Posto isto, é possível criar dois vetores normais à linha, simétricos e com a mesma intensidade, cuja soma das intensidades dos dois vetores equivale à margem – estes vetores são denominados de **vetores de suporte** (em inglês, *support vector*) [27].

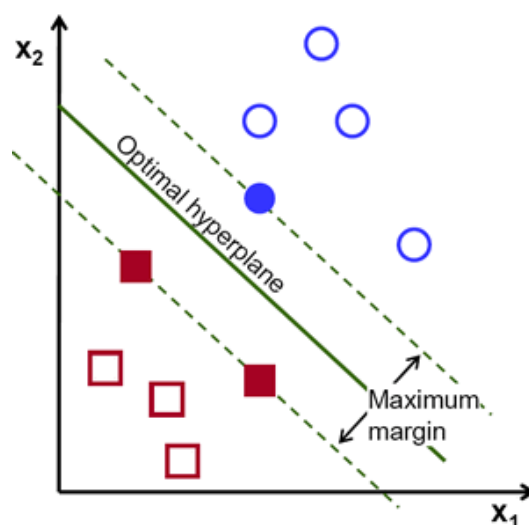


Figura 9: Exemplo de uma SVM. Fonte: [28]

Quando as observações estão distribuídas de uma tal forma que não é possível obter um hiperplano capaz de as separar (isto é, as observações não são linearmente separáveis), existe a necessidade de mudar o domínio do problema de tal forma que no novo domínio seja possível existir o hiperplano – tal como se encontra representado na Figura 10. A esta mudança de domínio para um domínio superior, esta transformação que se aplica é denominada de **truque do núcleo** (em inglês, *kernel trick* ou *kernel method*) ou **nuclearização** (em inglês, *kerneling*).

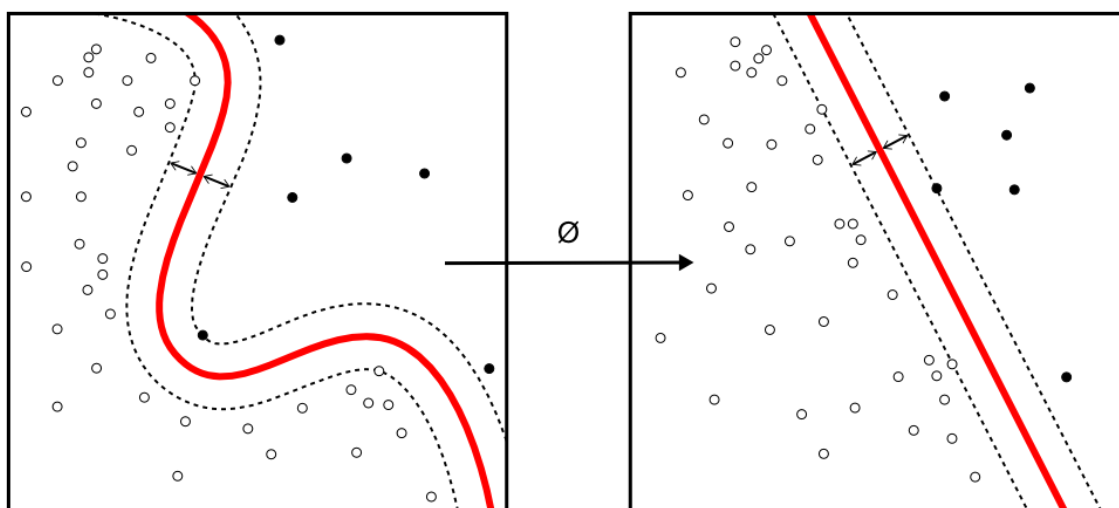


Figura 10: Aplicação de uma função de nuclearização. Fonte: [29]

Se as observações forem classificadas por mais de duas classes, tipicamente divide-se o problema de classificação em vários sub-problemas de classificação através do emparelhamento de classes distintas em cada problema. Os emparelhamentos mais conhecidos são os **Um-Contra-Um** (em inglês, *One-versus-One*) e **Um-Contra-Todos** (em inglês, *One-versus-All*): o primeiro emparelhamento tem como fundamento seleccionar o sub-classificador com a maior soma de positivos previstos, enquanto que o segundo emparelhamento tem como fundamento seleccionar o sub-classificador com a maior pontuação (entre 0 e 1) prevista.

2.3.4 Floresta Aleatória

Uma **floresta aleatória** (em inglês, *random forest*; RF como acrónimo) é um algoritmo de aprendizagem supervisionada, que pode ser considerada tanto de classificação como de regressão.

As florestas aleatórias consideram-se uma extensão das Árvores de Decisão, dado o problema que as árvores geram de sobreposição (em inglês, *overfitting*) dos dados de treino; isto é, quando uma árvore alcança uma profundidade considerável, é sinal que aprendeu padrões irregulares – o que vai tornar a variância das previsões bastante alta [30][31][32].

Para combater a sobreposição, gera-se um conjunto de árvores de decisão em que o conjunto de observações e as suas classificações de cada árvore é obtido selecionando aleatoriamente dados do conjunto fornecido de observações X e as suas classificações Y.

Dada uma nova observação, calcula-se a classificação prevista para cada árvore da floresta. A classificação mais prevista em toda a floresta é aquela que a floresta classificará a observação.

A Figura 11 ilustra um exemplo para uma floresta aleatória.

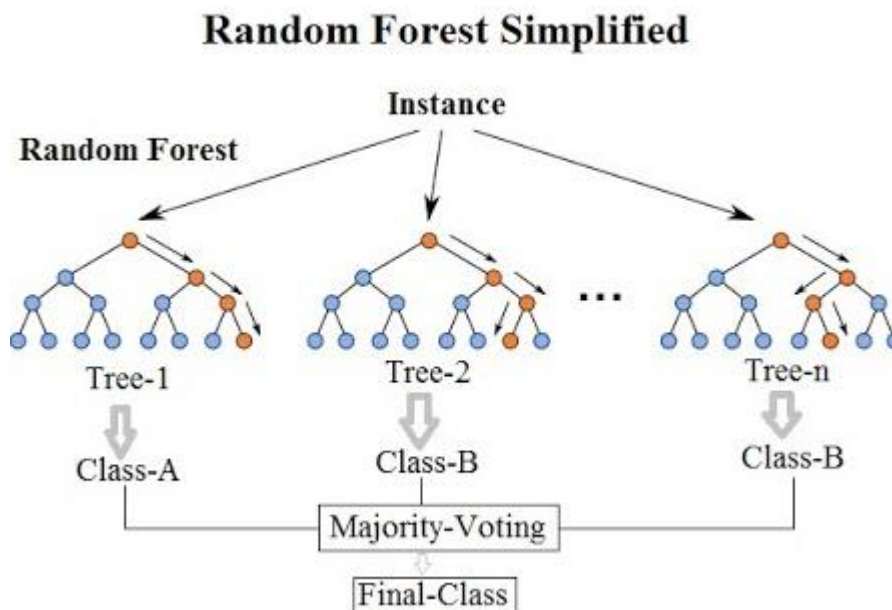


Figura 11: Exemplo de uma floresta aleatória. Fonte: [32]

2.3.5 Naive Bayes

Naive Bayes (NB como acrónimo) é um algoritmo de aprendizagem supervisionada, nomeadamente de classificação [33][35].

Este algoritmo faz parte de uma família de classificadores probabilísticos que se baseiam no **teorema de Bayes** [34][36]. Este teorema, cuja fórmula se encontra na Equação 5, indica a probabilidade de um evento A acontecer após ter ocorrido um evento B. Esta probabilidade é obtida através da divisão entre a probabilidade de ambos os eventos ocorrerem e a probabilidade do evento B (o primeiro evento a ocorrer).

$$p(A | B) = \frac{p(A \cap B)}{p(B)} = \frac{p(B | A) p(A)}{p(B)}$$

Equação 5: Teorema de Bayes.

Utilizando a terminologia da probabilidade bayesiana, a Equação 5 pode ser simplificada na maneira como se encontra ilustrado na Equação 6.

$$posteriori = \frac{observado * priori}{evidencia}$$

Equação 6: Teorema de Bayes utilizando a terminologia da probabilidade bayesiana.

Relacionado o teorema com o conjunto de observações X e as suas classificações Y (representado na Equação 7): o objetivo do algoritmo é calcular a probabilidade de uma classificação em Y ser um valor Y_v quando os atributos X_i da observação são visíveis.

$$p(Y_v | X_i) = \frac{p(X_i | Y_v) p(Y_v)}{p(X_i)}$$

Equação 7: Teorema de Bayes aplicado ao problema de classificação.

E ainda a fórmula da Equação 7 pode ser simplificada segundo a probabilidade em cadeia (como é possível visualizar na Equação 8), assumindo que os atributos são condicionalmente independentes.

$$p(Y_v | X_i) = p(X_1 | Y_v) * p(X_2 | Y_v) * ... * p(X_n | Y_v), X_i = \{X_1, X_2, ..., X_n\}$$

Equação 8: Teorema de Bayes e a probabilidade em cadeia.

Posto isto, para construir o classificador há necessidade de realizar três passos. O primeiro passo é construir uma **tabela de frequência** (como o exemplo da Tabela 1) para cada atributo A das observações X e respetivo valor, relacionando com a classificação final; isto é, contar o número de classificações distintas para cada par atributo-valor.

Tabela de Frequência		Joga?	
		Sim	Não
Estado do Tempo	Sol	3	2
	Nublado	4	0
	Chuva	2	3

Tabela 1: Exemplo de uma tabela de frequência para o atributo Estado do Tempo.

Fonte: [36]

De seguida, construir uma **tabela de probabilidades** para cada atributo A , calculando as probabilidades à posteriori para cada valor do atributo A e as classificações possíveis. No exemplo da Tabela 2, as probabilidades a azul indicam a probabilidades de evidência; as probabilidades a verde indicam as probabilidades à priori e as restantes probabilidades indicam as probabilidades já observáveis.

Tabela de Probabilidades		Joga?		
		Sim	Não	
Estado do Tempo	Sol	3/9	2/5	5/14
	Nublado	4/9	0/5	4/14
	Chuva	2/9	3/5	5/14
		9/14	5/14	

Tabela 2: Exemplo de uma tabela de probabilidades para o atributo Estado do Tempo. Fonte: [36]

Existem duas exceções à regra: tipicamente utiliza-se a distribuição Gaussiana para construir a tabela de probabilidades para atributos contínuos e utiliza-se a distribuição Bernoulli para construir a tabela de probabilidades para atributos booleanos.

Uma vez construídas as tabelas dos passos anteriores e calculadas as probabilidades à posteriori de todos atributos, o classificador encontra-se construído. O objetivo agora é dado uma observação, classifica-la. Isto é possível calculando a probabilidade à posteriori para cada classificação, utilizando a equação da regra em cadeia para os respectivos valores dos atributos fornecidos. Uma vez as probabilidades normalizadas (isto é, cada probabilidade é dividida pela soma de todas as probabilidades à posteriori calculadas), a probabilidade com maior valor é aquela que classificará a observação fornecida.

2.3.6 K-Means

K-Means é um algoritmo de aprendizagem não-supervisionada, cujo foco é o agrupamento. Isto é, dado um conjunto de observações X em que cada observação é composta por N atributos, o objetivo é agrupar as observações em K categorias ($K < N$) [37][38].

Para tal, é fornecido um número K de observações (definidos ou aleatórios), conhecidos como **centroides** (em inglês, *centroids*), de forma que o algoritmo modifique os valores/posições dos centroides capaz de minimizar a variância.

Por cada iteração, o algoritmo adiciona uma observação X_i a um conjunto de observações X'_j mais próximas de cada centroide C_j , utilizando uma métrica D (como distância euclidiana ou distância de Manhattan) como se encontra definido na Equação 9.

$$\arg \min_j D(X_i, C_j)$$

Equação 9: Fórmula para obter observações mais próximas de cada centroide.

Ao definir todos os conjuntos de observações, o algoritmo modifica os centroides com base nos valores/posições médias de cada conjunto, tal como se encontra definido na Equação 10. O algoritmo termina quando os centroides convergirem de tal forma que pouco se vão

modificar com mais iterações. As categorias encontram-se definidas com base nos valores dos centroides.

$$C_j = \frac{1}{N_j} \sum_{k=1}^{N_j} (X'_{j_k})$$

Equação 10: Fórmula para modificar cada centroide.

2.3.7 GMM

GMM (acrónimo para *Gaussian Mixture Model*; em português, Modelo de Mistura Gaussiano) é um algoritmo de aprendizagem não-supervisionada, cujo foco é o agrupamento [39][40][41][42][43].

O objetivo deste algoritmo é moldar a distribuição das observações no espaço de entrada, organizando os dados de acordo com a componente a que fiquem associados.

Apenas considerando que as observações são compostas por um único atributo, a Equação 11 ilustra como calcular a probabilidade dessa observação, sabendo que existem K componentes e que a cada componente está associado um peso Φ e uma distribuição normal $N(\mu, \sigma^2)$, cuja fórmula de cálculo da probabilidade nessa distribuição se encontra na Equação 12.

$$p(x) = \sum_{i=1}^K (\Phi_i * N(x | \mu_i, \sigma_i^2))$$

Equação 11: Fórmula da probabilidade de um GMM.

$$N(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} * e^{-\frac{1}{2}\left(\frac{(x-\mu)^2}{\sigma^2}\right)}$$

Equação 12: Fórmula para a função de distribuição gaussiana.

No caso em que as observações X são compostas por D atributos, a Equação 13 ilustra como calcular a probabilidade de uma observação x (segundo a forma de um vetor de dimensão D), sendo que Equação 14 mostra como se obtém a probabilidade numa distribuição gaussiana para múltiplos atributos, sendo que S_i representa a matriz da covariância.

$$p(\vec{x}) = \sum_{n=1}^K (\Phi_n * N(\vec{x} | \vec{\mu}_n, S_n))$$

Equação 13: Fórmula da probabilidade de um GMM para múltiplos atributos.

$$N(\vec{x} | \vec{\mu}, S) = \frac{1}{\sqrt{(2\pi)^k |S|}} * e^{-\frac{1}{2}(\vec{x} - \vec{\mu})^T S^{-1} (\vec{x} - \vec{\mu})}$$

Equação 14: Fórmula para a função de distribuição gaussiana para múltiplos atributos.

O problema que se coloca é como obter os pesos Φ , as médias μ e os desvios padrões σ . Para tal, recorre-se a um algoritmo que estima quais os parâmetros com base nos dados de entrada – o algoritmo *Expectation-Maximization* (EM como acrónimo; em português, Espetativa-Maximização) [44].

Primeiro, inicializa-se para cada um dos K componentes as três variáveis a estimar. As médias μ são inicializadas escolhendo aleatoriamente um elemento dos dados de entrada X. De seguida, inicializa-se as variâncias σ^2 segundo a fórmula da Equação 15, em que \bar{x} é a média dos dados de entrada X. Por fim, inicializa-se os pesos Φ dividindo um pelo número K de componentes.

$$\sigma_i^2 = \frac{1}{N} \sum_{j=1}^N [(x_j - \bar{x})^2]$$

Equação 15: Fórmula para inicializar as variâncias.

Uma vez inicializadas as variáveis, dá-se início a várias iterações para modificar estas variáveis em que o algoritmo termina quando um certo limite de erro não for encontrado em todas as variáveis entre a última iteração coma penúltima iteração.

Por cada ciclo, calcula-se a probabilidade γ de um elemento x_i pertencer a uma componente k das K componentes, através da Equação 16.

$$\gamma_{ik} = p(C_k | x_i) = \frac{\Phi_i * N(x_i | \mu_k, \sigma_k^2)}{\sum_j^K (\Phi_j * N(x_i | \mu_j, \sigma_j^2))}$$

Equação 16: Fórmula da probabilidade de um elemento i pertencer a uma componente k.

Uma vez calculada, atualiza-se as três variáveis segundo as Equação 17, Equação 18 e Equação 19.

$$\Phi_k = \sum_{i=1}^N \left[\frac{\gamma_{ik}}{N} \right]$$

Equação 17: Fórmula para atualizar os pesos da fórmula do algoritmo GMM.

$$\mu_k = \frac{\sum_{i=1}^N [\gamma_{ik} * x_i]}{\sum_{i=1}^N [\gamma_{ik}]}$$

Equação 18: Fórmula para atualizar as médias da fórmula do algoritmo GMM.

$$\sigma_k = \frac{\sum_{i=1}^N [(x_i - \mu_i)^2]}{\sum_{i=1}^N [\gamma_{ik}]}$$

Equação 19: Fórmula para atualizar os desvios-padrões da fórmula do algoritmo GMM.

Uma vez concluído, é possível prever a que componente pertence uma nova observação x , calculando simplesmente a probabilidade de pertença em todas as componentes (ilustrado na Equação 16, sob a forma de $p(C_k | x)$) e escolher a probabilidade com maior valor.

2.4 Processamento de Imagem

Processamento de imagem (em inglês, *image processing*) consiste na análise e tratamento de imagem, aplicando uma diversidade de algoritmos matemáticos para o conseguir. Nesta seção, serão abordados os modelos de cor que representam as imagens, os algoritmos de segmentação de imagem e as características que serão extraídas, tanto da máscara obtida na segmentação tanto da imagem real [45].

2.4.1 Modelos de Cor

Um modelo de cor (em inglês britânico, *colour model*; em inglês americano, *color model*) corresponde ao sistema de cores utilizado para representar uma imagem. Tipicamente os modelos mais utilizados são o **modelo RGB** (acrónimo para *Red Green Blue*; em português, Vermelho Verde Azul) e o **modelo HSV** (acrónimo para *Hue Saturation Value*; em português, Matiz Saturação Valor) [46][47].

No modelo RGB, cada pixel encontra-se definido em três canais, correspondentes às três cores definidas acima. Cada canal tem um valor entre uma intensidade nula e uma intensidade máxima (tipicamente 255, devido ao mínimo ocupado por um pixel ser 1 Byte).

No modelo HSV, cada pixel encontra-se definido em três canais, correspondentes aos três parâmetros definidos acima. Cada canal tem um valor entre 0 e 1. A Matiz corresponde ao nível de cinzentos de uma cor, enquanto que a Saturação corresponde ao eixo entre a cor da Matiz pura e a cor selecionada do Valor puro.

2.4.2 Segmentação

Segmentação (em inglês, *segmentation*) consiste no processo de simplificação de imagem, através da partição da mesma em várias categorias ou grupos, com o intuito de obter uma visualização, uma versão da imagem mais simplificada e com melhor significado ou qualidade.

Este processo permite, sobretudo, obter os contornos, os limites e o conteúdo de objetos representados na imagem – isto é, a máscara (em inglês, *mask*) do objeto. Para tal, vários métodos de segmentação podem ser aplicados, recorrendo até a técnicas de pré-processamento da imagem, como **suavização** (em inglês, *smoothing*) ou **transformações geométricas**.

Nas seguintes subsecções, serão abordados vários algoritmos de segmentação, onde o foco será no processo de **binarização** (em inglês, *binarization*) da imagem: atribuição de pixéis a branco como o objeto identificado em **primeiro plano** (em inglês, *foreground*) e a atribuição de pixéis a preto com o resto da imagem – também conhecido como **fundo** (em inglês, *background*).

2.4.2.1 Método de *Thresholding*

O método de *Thresholding* tem como objetivo atribuir a uma imagem, numa escala de cinzentos no modelo RGB, a cor branca (intensidade máxima) a todos os pixéis com intensidade superior a uma intensidade limite θ , enquanto que serão atribuídos a cor preta (intensidade nula) a todos os pixéis com intensidade inferior ou igual a esse mesmo limite [48].

Este método pode ser considerado uma extensão da função de *thresholding* (ilustrada na Figura 12), em que os dados de entrada são os pixéis da imagem e os dados de saída são os mesmos pixéis com a sua intensidade alterada.

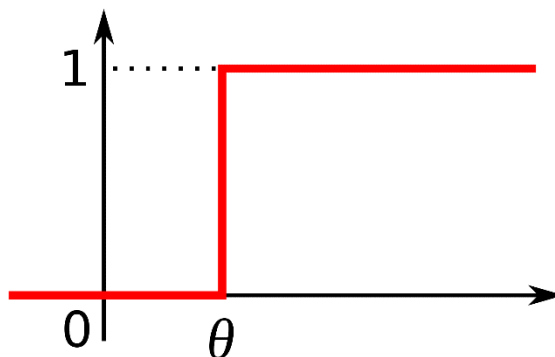


Figura 12: Função de *thresholding*. Fonte: [49]

2.4.2.2 Método de Otsu

O método de Otsu (em inglês, *Otsu's method*) tem como objetivo obter um valor limite (em inglês, *threshold*) de forma automática, para o aplicar no método de *thresholding* sobre uma dada imagem [50]. O método obtém automaticamente esse limite θ analisando o histograma de intensidades de uma imagem, encontrando a intensidade ideal entre dois picos de intensidades.

Percorrendo todas as intensidades t (entre 0 e 255, no caso do modelo RGB) da imagem guarda-se dois conjuntos: o conjunto 1 que contém todos os pixéis com intensidade inferior ou igual a t e o conjunto 2 que contém todos os pixéis com intensidade superior a t . Pretende-se

encontrar a intensidade M que minimiza a soma pesada entre a soma S_c de todas as intensidades com a variância σ_c^2 dos dois conjuntos. Esta minimização encontra-se definida na Equação 20.

$$\min f(t) = S_1(t) * \sigma_1^2(t) + S_2(t) * \sigma_2^2(t)$$

Equação 20: Fórmula da obtenção do *threshold*. Fonte: [51]

O valor t da minimização é o valor a aplicar no método de *thresholding*.

2.4.2.3 K-Means

Remodelando (em inglês, *reshaping*) os pixéis de uma imagem (numa escala de cinzentos) num só vetor, é possível aplicar o algoritmo de aprendizagem não-supervisionada K-Means para segmentar a imagem, tendo como resultado centroides que equivalem a uma intensidade. Quando o valor K do algoritmo é igual a 2, aplica-se o método de *thresholding* à intensidade que corresponde ao maior centroide menos um.

2.4.2.4 GMM

Tal e qual como no K-Means, remodelando os pixéis de uma imagem (numa escala de cinzentos) num só vetor, é possível aplicar o algoritmo de aprendizagem não-supervisionada GMM para segmentar a imagem, tendo como resultado componentes que equivalem a uma intensidade. Quando o valor K do algoritmo é igual a 2, aplica-se o método de *thresholding* à intensidade que corresponde à maior componente menos um.

2.4.2.5 Watershed

O algoritmo **Watershed** (em português, algo semelhante como linha divisória de água) tem como fundamento as marcas de água [52][53]; isto é, dado uma intensidade t e um limite de crescimento C , o algoritmo irá obter, primeiro, a imagem segmentada através do método de *thresholding* – para um limite igual a t . A imagem obtida também é conhecida por imagem raiz. De seguida, percorrerá C intensidades com valor superior ou inferior a t . Por cada iteração ou intensidade, obtém-se uma imagem segmentada através do mesmo método de *thresholding*.

Se os pixéis a branco da nova imagem segmentada forem adjacentes (segundo uma vizinhança definida) aos pixéis a branco da imagem raiz, então esses mesmos pixéis são incluídos na imagem raiz. Este processo pode ser exemplificado segundo a Figura 13.

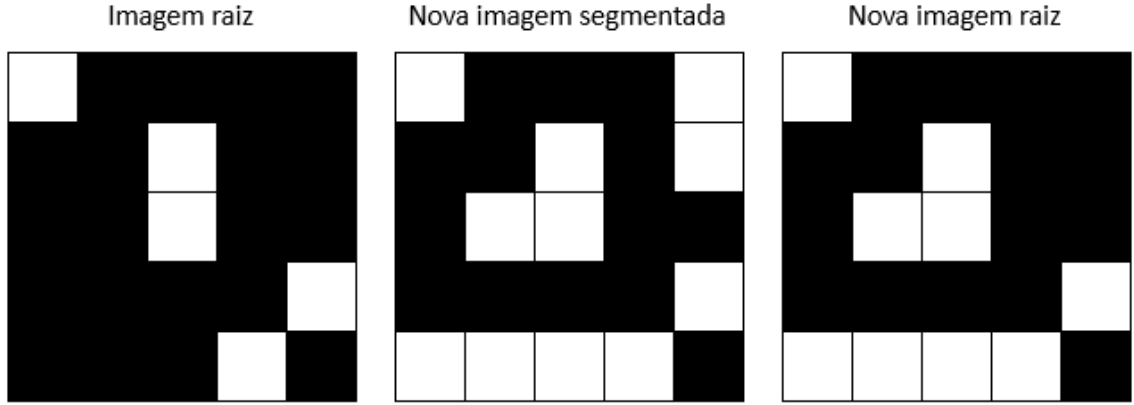


Figura 13: Preenchimento da imagem raiz pelo algoritmo de Watershed.

2.4.2.6 Cortes Normalizados

O método dos **cortes normalizados** (em inglês, *normalized cuts*; NCut como acrónimo) baseia-se no princípio de separar um grafo em dois, aplicando um corte mínimo. A ideia é olhar uma imagem como um grafo não-dirigido, em que os pixels representam vértices e cada ligação entre pixels equivale a uma aresta com um peso unitário.

Um grafo G , que representa uma imagem, é definido por um conjunto de vértices V e um conjunto de arestas E – pode ser igualmente definido na forma $G = (V, E)$. O objetivo é separar o conjunto de vértices V em dois conjuntos disjuntos de vértices A e B . Isto terá como consequência a criação de dois grafos G_1 e G_2 , em que G_1 é composto por A vértices e E_1 arestas e G_2 é composto por B vértices e E_2 arestas.

Um corte é definido como uma soma dos pesos das arestas entre os vértices A e os vértices B . Este valor é calculado segundo a Equação 21. O resultado suposto é encontrar o corte mínimo para encontrar obter a segmentação de um objeto. O problema que se coloca é como definir os vértices A e B e os pesos das ligações.

$$cut(A, B) = W(A, B) = \sum_{p \in A, q \in B} w_{pq}$$

Equação 21: Fórmula para o corte em A e B .

Uma solução para definir os vértices A e B é utilizar algoritmos que minimizam o corte da Equação 21, tais como o **algoritmo de Karger** ou o **algoritmo de Stoer-Wagner** [55][56].

Uma solução para a definição dos pesos das ligações é assumir que o peso mede a similaridade, obtida pela Equação 22, com base na função da distância D (que pode ser tanto euclidiana ou de Manhattan) e um parâmetro σ (que assume um valor entre 0 e 1) de fator de distância.

$$w_{pq} = e^{-\frac{1D(p,q)^2}{2\sigma^2}}$$

Equação 22: Fórmula para definir os pesos das ligações do grafo. Fonte: [58]

O problema que aqui se coloca é que os algoritmos que minimizam a Equação 21 tendem encontrar pequenos conjuntos de vértices para realizar a separação. Introduce-se então uma medida de dissociação denominada **corte normalizada** [57]. Esta medida pode ser calculada segundo a Equação 23.

$$Ncut(A, B) = \frac{W(A, B)}{W(A, V)} + \frac{W(A, B)}{W(B, V)}$$

Equação 23: Fórmula para calcular a medida de dissociação.

Da mesma maneira que se pode definir uma medida de dissociação, também se pode definir uma medida de associação, que pode ser definida na Equação 24. Esta relação entre medidas pode ser visível na Equação 25.

$$Nassoc(A, B) = \frac{W(A, A)}{W(A, V)} + \frac{W(B, B)}{W(B, V)}$$

Equação 24: Fórmula para calcular a medida de associação.

$$Ncut(A, B) + Nassoc(A, B) = 2$$

Equação 25: Relação entre a medida de associação e a medida de dissociação.

Dada a introdução destes dois conceitos, resta saber quais os vértices A e B; ou seja, em que arestas do grafo se vai realizar o corte.

Primeiro, define-se a matriz de adjacência W – uma matriz quadrada simétrica de dimensão igual a |V| e um elemento (i, j) da matriz corresponde ao número de arestas que existe na direção de i para j. De seguida, define-se uma matriz diagonal D – uma matriz quadrada em que todos os elementos que não pertençam à diagonal principal da matriz são iguais a zero. Esta diagonal é definida segundo a Equação 26.

$$D_{ii} = \sum_{j=1} W_{ij}$$

Equação 26: Definição da matriz diagonal.

Uma vez definida a matriz D, define-se um vetor indicador y com um tamanho igual ao número de vértices do grafo G, em que o elemento y_i tem valor 1 se o vértice i pertencer a A e valor negativo (tipicamente, -1) se pertencer a B. Posto isto, substituindo na Equação 21 por estes novos dados, o objetivo é minimizar o valor do corte definido na Equação 27.

$$cut(A, B) = W(A, B) = \frac{y^T(D - W)y}{y^T D y}$$

Equação 27: Nova fórmula do corte para o algoritmo do corte normalizado.

Um novo problema se encontra minimizando o valor do corte: o algoritmo proposto de minimização é do tipo *NP-Hard*. Para combater esta situação, se o vetor y for relaxado para assumir valores reais, calculam-se os **vetores próprios** (em inglês, *eigenvector*), segundo a Equação 28, para calcular o segundo menor **valor próprio** (em inglês, *eigenvalue*).

$$(D - W)y = \lambda Dy$$

Equação 28: Equação para minimizar o corte normalizado, calculando os vetores próprios.

Este valor próprio calculado será utilizado para bipartir o grafo e, conseqüentemente, a imagem. Caso haja necessidade, o mesmo algoritmo pode ser utilizado para um novo grafo e assim sucessivamente.

2.4.2.7 Closed Shortest Path

O método *Closed Shortest Path* (em português, caminho fechado mais curto) baseia-se nas técnicas de pesquisa do caminho mais curto aplicados a grafos. A ideia é olhar uma imagem como um grafo não-dirigido, em que os pixéis representam vértices e cada ligação entre pixéis equivale a uma aresta com um peso unitário.

Um grafo não-dirigido G , que representa uma imagem, é definido por um conjunto de vértices V e um conjunto de arestas E – pode ser igualmente definido na forma $G = (V, E)$. Um **grafo acíclico dirigido** (em inglês, *directed acyclic graph*; DAG como acrónimo) é um grafo dirigido em que não se verifica ligações entre arestas capazes de formar um ciclo. Tipicamente, transforma-se um grafo não-dirigido numa DAG com recurso à **ordenação topológica** (em inglês, *topological sorting*), utilizando algoritmos como **pesquisa por profundidade** (em inglês, *deep-first search*; DFS como acrónimo) ou o **algoritmo de Kahn** [59][60].

A vizinhança (em inglês, *neighborhood*) de um vértice são todos os vértices a uma dada distância mínima desse mesmo vértice, aquilo que se chama de conectividade. Os vizinhos causais (em inglês, *causal neighbors*) de um vértice são todos os vizinhos que antecedem o vértice.

O método tem como objetivo definir um caminho com o menor custo à volta de um **ponto de raiz** (em inglês, *seed point*) definido, com base em minimizar uma função de soma de medidas como a **magnitude da derivada** (em inglês, *magnitude of derivate*) ou a **distância euclidiana**. Nesta versão do método, aplica-se as coordenadas originais e não as coordenadas polares [61].

Primeiro, construi-se o grafo não-dirigido G com base nos pixéis da imagem em questão com os pesos das arestas a assumir um valor que é obtido pela Equação 29, com base na magnitude da derivada d no vértice de entrada (ou cabeça da aresta) e com base em três constantes: a , b e c .

$$w = a + (b - a) * \frac{e^{(255-d)c} - 1}{e^{255c} - 1}$$

Equação 29: Definição do peso de uma aresta que parte de um vértice.

De seguida, define-se o ponto de raiz R (tipicamente ao centro da imagem) e realiza-se a linearização dos pixéis para uma dada vizinhança de raio r.

A Figura 14 exemplifica como seria a ordem dos pixéis a visitar, para uma vizinhança de raio 5. Todos os pixéis com valor nulo representam pixéis que se encontram fora do raio, exceto o pixel ao centro que representa o ponto de raiz R.

Se se considerar a vizinhança de conectividade 8, a vizinhança causal do vértice na posição 3 é composta pelo vértice na posição 2; a vizinhança causal do vértice na posição 7 é composta pelos vértices nas posições 2, 3, 4 e 6; a vizinhança causal do vértice na posição 79 é composta pelos vértices nas posições 73, 76, 77 e 78.

0	0	0	0	0	21	0	0	0	0	0
0	0	31	28	26	22	20	18	15	0	0
0	35	32	30	27	23	19	16	14	11	0
0	37	36	33	29	24	17	13	10	9	0
0	40	39	38	34	25	12	8	7	6	0
41	42	43	44	45	0	1	2	3	4	5
0	46	47	49	54	61	72	77	79	80	0
0	48	50	53	58	62	68	73	76	78	0
0	51	52	56	59	63	67	70	74	75	0
0	0	55	57	60	64	66	69	71	0	0
0	0	0	0	0	65	0	0	0	0	0

Figura 14: Exemplo da linearização com base numa vizinhança de raio 5. Fonte: [61]

Define-se uma DAG D com a ordem como os vértices estão ordenados na linearização. Para realizar o fecho do caminho, introduz-se de novos os pixéis quando o ângulo é igual a zero. Exemplificado com a Figura 14, após inserir normalmente os vértices adiciona-se os vértices nas posições 1, 2, 4 e 5 na respetiva vizinhança das posições 61, 72, 77, 79 e 80. Por exemplo, a vizinhança procedente do vértice na posição 77 são os vértices das posições 2, 3 e 79; todos os vértices das posições 2 e 3, nesta vizinhança, são tomados como fim de ciclo.

Agora é calcular o caminho mais curto para a DAG D com base no peso das arestas definidas. O resultado será um conjunto de vértices que correspondem a pixéis de um contorno. Todo o conteúdo dentro do contorno tomará intensidade máxima enquanto que o resto da imagem tomará intensidade nula.

2.4.3 Extração de Características

Extração de Características (em inglês, *feature extraction*) consiste no processo em obter informação específica de uma imagem, informação essa que a permita caracterizar. Cada problema em que um conjunto de imagens estão inseridas têm uma certa quantidade de características a extrair.

Aquando estudar objetos segmentados de imagens, na Tabela 3 estão indicadas as características a extrair dos objetos segmentados e em que categorias pertencem [62].

Categoria Principal	Subcategoria	Características	Quantidade
Propriedades do Contorno	Propriedades Geométricos	Área	1
		Perímetro	1
		Diâmetro maior	1
		Diâmetro menor	1
		Diâmetro Equivalente	1
		Compacidade	1
		Circularidade	1
		Solidez	1
		Retangularidade	1
		Proporção	1
		Excentricidade	1
	Assimetria da Segmentação	Rácio Médio	1
		Desvio-Padrão do Rácio	1
Propriedades da Cor	Modelo de Cor	Cor Média	3
		Desvio-Padrão da Cor	3
		Intensidade Mínima	3
		Intensidade Máxima	3
		Inclinação da Cor	3
Análise da Textura	Matriz de Coocorrências	ASM	12
		Correlação	12
		Contraste	12

Tabela 3: Categorias de características relacionadas com objetos segmentados.

Nas seguintes subseções serão abordadas as várias características em cada subcategoria.

2.4.3.1 Propriedades Geométricas

A **área** (em inglês, *area*) de um contorno é uma característica que é obtida contando todos os pixels dentro de um contorno; o **perímetro** (em inglês, *perimeter*) igualmente é uma característica que é obtida contando todos os pontos do contorno.

O **diâmetro maior** (em inglês, *major diameter*) é a maior distância entre dois pontos do contorno, enquanto que o **diâmetro menor** (em inglês, *minor diameter*) é a menor distância entre dois pontos do contorno.

O **diâmetro equivalente** (em inglês, *equivalent diameter*) corresponde ao diâmetro de um círculo que contem a mesma área que o contorno. A Equação 30 mostra como se calcula esta característica, baseando apenas na área A do perímetro.

$$DE = \sqrt{\frac{4A}{\pi}}$$

Equação 30: Fórmula do diâmetro equivalente.

A **compacidade** (em inglês, *compactness*) corresponde à proporção entre a área do contorno e o perímetro de um círculo com a mesma área do contorno. A Equação 31 mostra como se calcula esta característica, baseando apenas no diâmetro equivalente DE e no maior diâmetro MD.

$$CO = \frac{DE}{MD}$$

Equação 31: Fórmula da compacidade.

A **circularidade** (em inglês, *circularity*) corresponde ao quanto o contorno se aproxima de um círculo. A Equação 32 mostra como se calcula esta característica, baseando apenas na área A e no perímetro P.

$$CI = \frac{4\pi A}{P^2}$$

Equação 32: Fórmula da circularidade.

A **solidez** (em inglês, *solidity*) corresponde à proporção entre a área do contorno e a área obtida após tornar o contorno convexo. A Equação 33 mostra como se calcula esta característica, baseando apenas na área A e na área A_C do contorno convexo.

$$SO = \frac{A}{A_C}$$

Equação 33: Fórmula da solidez.

A **retangularidade** (em inglês, *extent*) corresponde à proporção entre a área do contorno e a área de um retângulo que cobre o contorno. A Equação 34 mostra como se calcula esta característica, baseando apenas na área A e na área A_R do retângulo que cobre o contorno.

$$SO = \frac{A}{A_R}$$

Equação 34: Fórmula da retangularidade.

A proporção (em inglês, *aspect ratio*) corresponde à divisão entre o maior diâmetro e o menor diâmetro. A Equação 35 mostra como se calcula esta característica, baseando apenas no maior diâmetro MD e no menor diâmetro mD.

$$PR = \frac{MD}{mD}$$

Equação 35: Fórmula da proporção.

A excentricidade (em inglês, *eccentricity*) corresponde à medida da elongação da região envolvente do contorno. A Equação 36 mostra como se calcula esta característica, baseando apenas nos momentos Mu calculados do contorno.

$$e = \frac{4Mu_{11} * (Mu_{02} - Mu_{20})^2}{(Mu_{02} + Mu_{20})^2}$$

Equação 36: Fórmula da excentricidade.

O momento Mu_{ij} pode ser calculado segundo a Equação 37, sendo que a função f corresponde à intensidade do pixel na posição (x,y) [65].

$$Mu_{ij} = \sum_x \sum_y x^i y^j f(x, y)$$

Equação 37: Fórmula do momento ij.

2.4.3.2 Assimetria de um objeto

A **assimetria de um objeto** (em inglês, *object asymmetry*) corresponde à ausência de simetria; isto é, à ausência de linhas capazes de separar uma imagem em duas imagens e afirmar que uma é espelho da outra.

Uma forma típica de estimar a assimetria consiste em selecionar os pontos desse diâmetro (que estão dentro do contorno) e calcula-se os eixos perpendiculares ao maior diâmetro nesses mesmos pontos. Calcula-se a **proporção** entre o segmento de reta acima do maior diâmetro e o segmento de reta abaixo do maior diâmetro, para cada ponto do maior diâmetro.

A assimetria de um objeto fica definida com a **média** (em inglês, *average* ou *mean*) e o **desvio-padrão** (em inglês, *standard deviation*) das proporções calculadas. A Figura 15 ilustra como são calculados os segmentos de reta.

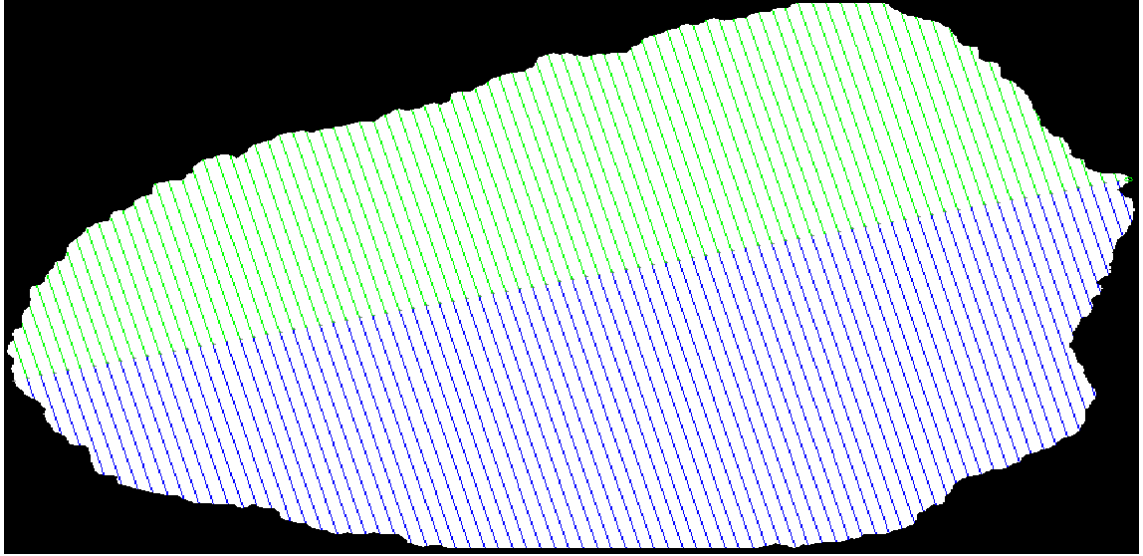


Figura 15: Exemplo da obtenção dos segmentos de reta.

2.4.3.3 Propriedades da Cor

Quanto às propriedades da cor, pretende-se calcular características para um dado modelo de cor. Independentemente do modelo, é importante obter, como características, a **cor média** dos pixels dentro do contorno e o respetivo **desvio-padrão**, a **intensidade máxima** e **intensidade mínima** por cada canal de cor e a **inclinação** (em inglês, *skewness*) da cor para cada canal.

A inclinação de um conjunto de dados encontra-se definida na Equação 38, baseando apenas num número N de dados, na média μ e no desvio-padrão σ .

$$INC = \frac{\sum_{p=1}^N (I_p - \mu)^3}{N\sigma^3}$$

Equação 38: Fórmula da Inclinação.

2.4.3.4 Matriz de Coocorrências

Uma **matriz de coocorrências** (em inglês, *Co-Occurrence Matrix*) é uma matriz quadrada de ordem igual ao número de intensidades que um canal de cor consegue armazenar [63][64]. O objetivo desta matriz é determinar, para uma dada direção D (tipicamente: 0°, 45°, 90° e 135°), quantas vezes uma intensidade A é seguida da intensidade B. Na Figura 16 encontra-se dois exemplos de como foram calculadas duas matrizes de coocorrências com direção 0° (ou seja, da esquerda para a direita).

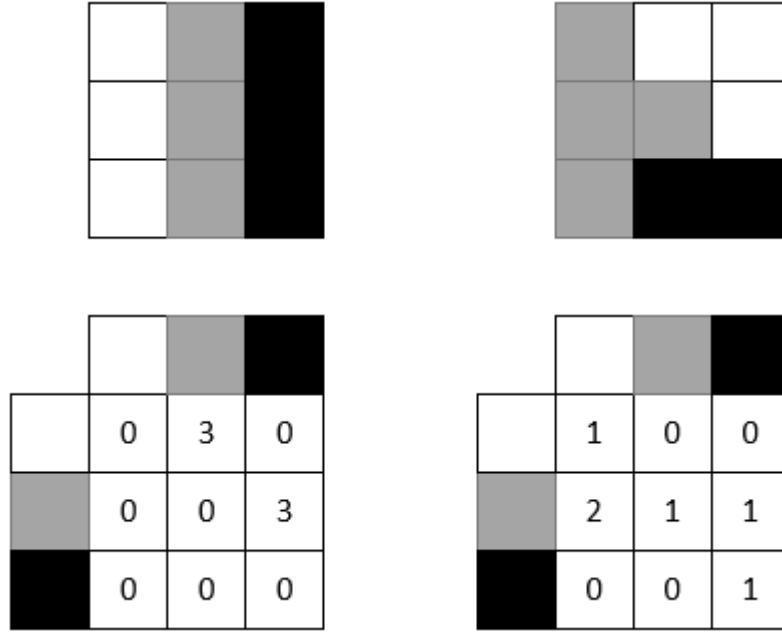


Figura 16: Exemplo de duas matrizes de ocorrência geradas a partir de duas imagens.

As características a extrair das matrizes de coocorrência são o **segundo momento angular** (em inglês, *angular second moment*), o **contraste** (em inglês, *contrast*) e a **correlação** (em inglês, *correlation*).

Ambas as características são obtidas a partir, respetivamente, das Equação 39, Equação 40 e Equação 41, em que N corresponde à ordem da matriz quadrada m , μ_i e σ_i corresponde à média e desvio-padrão dos elementos da linha i e μ_j e σ_j corresponde à média e desvio-padrão dos elementos da coluna j .

$$ASM = \sum_{i=1}^N \sum_{j=1}^N (m_{ij})^2$$

Equação 39: Fórmula do segundo momento angular.

$$C = \sum_{i=1}^N \sum_{j=1}^N m_{ij} (i - j)^2$$

Equação 40: Fórmula do contraste.

$$CRL = \sum_{i=1}^N \sum_{j=1}^N m_{ij} \left(\frac{(i - \mu_i)(j - \mu_j)}{\sqrt{\sigma_i^2 \sigma_j^2}} \right)$$

Equação 41: Fórmula da correlação.

2.5 Métricas

Métricas (em inglês, *metrics*) são medidas de desempenho dos vários algoritmos/técnicas utilizadas e estudadas anteriormente.

Uma das maneiras mais fácil de entender o conceito de métricas é, primeiro, construir uma **matriz de confusão** ou uma matriz de erro (em inglês, *confusion matrix*) para comparar os dados obtidos (ou previstos) de uma dada técnica com os dados reais de referência (em inglês, *ground truth data*).

A Tabela 4 representa o conceito de matriz de confusão [66]. Para um dado contexto, a **Condição Prevista** indica se a técnica em causa considerou o resultado obtido como verdadeiro (dando a indicação de Positivo) ou falso (dando a indicação de Negativo); a **Condição Real** indica se o resultado obtido pela técnica é realmente correto (dando a indicação de Verdadeiro) ou incorreto (dando a indicação de Falso).

		Condição Prevista	
		Positiva	Negativa
Condição Real	Positivo	Verdadeiro Positivo	Falso Negativo
	Negativo	Falso Positivo	Verdadeiro Negativo

Tabela 4: Matriz de Confusão.

Analisando a Tabela 4, pode-se afirmar que:

- os elementos **Verdadeiros Positivos** (em inglês, *True Positive*; VP como acrónimo) são todos aqueles que a técnica considerou positivo corretamente; também são conhecidos como **acertos**.
- os elementos **Verdadeiros Negativos** (em inglês, *True Negative*; VN como acrónimo) são todos aqueles que a técnica considerou negativo corretamente; também são conhecidos como **rejeições corretas**.
- os elementos **Falsos Positivos** (em inglês, *False Positive*; FP como acrónimo) são todos aqueles que a técnica considerou positivo incorretamente; também são conhecidos como **falsos alarmes** ou **erros do tipo 1**.
- os elementos **Falsos Negativos** (em inglês, *False Negative*; FN como acrónimo) são todos aqueles que a técnica considerou negativo incorretamente; também são conhecidos como **falhas** ou **erros do tipo 2**.
- a soma de todos os quatro tipos de elementos mencionados anteriormente tem que igualar o número de elementos total T em causa.

Assim sendo, várias métricas podem ser retiradas da matriz de confusão. Podem ser métricas relacionadas com a condição real, métricas relacionadas com a condição prevista e métricas relacionadas numa maneira geral.

Relacionadas com a condição real:

- a **taxa de verdadeiros positivos**, **taxa de acerto** ou **sensibilidade** (em inglês, *sensitivity* ou *recall*) é também conhecida como a probabilidade de deteção. A Equação 42 traduz esta métrica.

$$TVP = \frac{VP}{VP + FN}$$

Equação 42: Equação da Sensibilidade.

- a **taxa de verdadeiros negativos** ou **especificidade** (em inglês, *specificity*). A Equação 43 traduz esta métrica.

$$TVN = \frac{VN}{VN + FP}$$

Equação 43: Equação da Especificidade.

- a **taxa de falsos positivos** ou **queda** (em inglês, *fall-out*) é também conhecida como a probabilidade do falso alarme. A Equação 44 traduz esta métrica.

$$TFP = \frac{FP}{VN + FP}$$

Equação 44: Equação da Queda.

- a **taxa de falsos negativos** ou **taxa de falha** (em inglês, *miss rate*). A Equação 45 traduz esta métrica.

$$TFN = \frac{FN}{VP + FN}$$

Equação 45: Equação da Taxa de Falha.

Relacionadas com a condição prevista:

- o **valor de previsão positivo** ou **precisão** (em inglês, *precision*). A Equação 46 traduz esta métrica.

$$VPP = \frac{VP}{VP + FP}$$

Equação 46: Equação da Precisão.

- o **valor de previsão negativa**. A Equação 47 traduz esta métrica.

$$VPN = \frac{VN}{VN + FN}$$

Equação 47: Equação do Valor de Previsão Negativa.

- a **taxa de descobrimento falso**. A Equação 48 traduz esta métrica.

$$VDF = \frac{FP}{VP + FP}$$

Equação 48: Equação da Taxa de Descobrimento Falso.

- a **taxa de omissão falsa**. A Equação 49 traduz esta métrica.

$$VOP = \frac{FN}{VN + FN}$$

Equação 49: Equação da Taxa de Omissão Falsa.

Relacionadas de uma maneira geral:

- a **exatidão** (em inglês, *accuracy*). A Equação 50 traduz esta métrica.

$$E = \frac{VP + VN}{T}$$

Equação 50: Equação da Exatidão.

- a **prevalência** (em inglês, *prevalence*). A Equação 51 traduz esta métrica.

$$P = \frac{VP + FN}{T}$$

Equação 51: Equação da Prevalência.

- o **índice de Jaccard** (em inglês, *Jaccard index*) [67][68]. A Equação 52 traduz esta métrica.

$$JI = \frac{VP}{VP + FP + FN}$$

Equação 52: Equação do Índice de Jaccard.

- o **índice de Sørensen-Dice**, **coeficiente de Sørensen-Dice** ou **valor F1** (em inglês, *Sørensen-Dice index*) [69][70]. A Equação 53 traduz esta métrica.

$$SI = 2 \frac{TVP * VPP}{TVP + VPP} = \frac{2 VP}{2 VP + FP + FN}$$

Equação 53: Equação do Índice de Sørensen-Dice.

A aplicação da matriz de confusão, no contexto da segmentação, consiste em saber quantos pixels a branco existem na máscara obtida e que realmente identificam a lesão na imagem. No contexto da aprendizagem automática, consiste em saber em quantas imagens os vários algoritmos identificaram melanoma e que identificaram corretamente.

2.6 Trabalhos Relacionados

Nesta seção será abordada trabalhos relacionados com a dissertação; isto é, trabalhos que tenham o mesmo objetivo que esta dissertação e que implementaram as suas soluções ao problema.

Entre muitos artigos científicos que mostram como vários algoritmos de segmentação foram implementados (e com resultados), o artigo *Comparison of Segmentation Methods for Melanoma Diagnosis in Dermoscopy Images* foi o ponto de partida para esta dissertação, pois faz um estudo comparativo entre vários algoritmos de segmentação aplicados à segmentação da lesão onde está identificada melanoma [71].

A dissertação *Pattern Recognition in Pigmented Skin Lesion Images using Ensemble Methods* faz um estudo mais aprofundado das lesões da pele, tanto ao nível da segmentação como ao nível da extração de características e classificação. Apesar desta dissertação utilizar a mesma base de dados, houve uma seleção mais rigorosa de imagens, devido ao facto de que existem lesões cujo bordo se encontrava fora da imagem; assim, nesta dissertação utilizou-se 1104 das 2000 imagens disponíveis e 916 das 1104 imagens foram identificadas melanoma. Igualmente todas as imagens utilizadas foram redimensionadas para 400x299, para simplificar o processamento. Esta dissertação essencialmente detalhou e mostrou as características essenciais a extrair, apesar de certas características que foram utilizadas, como a área e o perímetro de um contorno, não serem independentes da escala a que o dermatoscópio obteve as imagens [72].

A tese de doutoramento *Automatic Detection of Melanomas Using Dermoscopy Images* igualmente detalhou significativamente os algoritmos de segmentação e das características a extrair. Esta tese de doutoramento enquadrou-se na edição 2016 do desafio da ISIC, sendo que

esta tese de doutoramento trabalhou sobre todas as 1279 imagens dermatoscópicas, sendo que apenas 248 das 1279 estão identificadas melanoma [73].

Por fim, a pesquisa ***Malignant Melanoma Detection Based on Machine Learning Techniques: A Survey*** permitiu reunir num só documento os algoritmos, características e métricas utilizadas e calculadas em vários trabalhos relacionados com o tema da dissertação, reforçando a ideia de se focar tanto na maioria dos algoritmos de segmentação e aprendizagem automática, mencionados nesta dissertação, como também reforçar a ideia de extrair características de certas categorias de características, igualmente mencionadas nesta dissertação [74].

Capítulo 3

Metodologia

3.1 Introdução

Neste capítulo, será abordada a especificação da dissertação; isto é, a análise do problema em que a dissertação está inserida, a definição e detalhes da solução proposta para o problema e a validação da solução.

3.2 Análise do Problema

Tal como referido na seção *Motivação e objetivos* do Capítulo 1, a construção do algoritmo para identificar melanoma automaticamente foi elaborado com recurso a imagens dermatoscópicas fornecidas pelo desafio *Skin Lesion Analysis Towards Melanoma Detection*.

Este desafio fornece, para cada uma das três fases, um conjunto de dados de entrada e um respetivo conjunto de dados reais de referência – também conhecido por informação clínica ou diagnóstico. Independentemente da fase, são fornecidas 2000 imagens dermatoscópicas e um ficheiro CSV com a informação do paciente a quem foi obtida a respetiva imagem dermatoscópica.

Em questões de validação, o desafio fornece igualmente um conjunto de dados para cada fase:

- na fase de segmentação, são fornecidas 2000 máscaras das lesões da pele consideradas de referência;
- na fase de extração de características, é fornecido um ficheiro CSV com as características de referência, no formato JSON (acrónimo para *JavaScript Object*

Notation; em português, Objeto com Notação JavaScript), para cada uma das 2000 imagens dermatoscópicas;

- na fase de classificação, é fornecido um ficheiro CSV que indica, para cada imagem, a existência de melanoma, de nevo ou de ceratose seborreica através de duas colunas B e C. Se a célula da coluna B for 0 e da coluna C for 0, foi identificado nevo na imagem; se a célula da coluna B for 1, foi identificado melanoma na imagem; se a célula da coluna C for 1, foi identificado ceratose seborreica; não existe uma imagem com estas duas colunas preenchidas a 1 simultaneamente.

3.3 Proposta de Solução

Nesta seção, será proposta e detalhada uma solução para o problema analisado na seção anterior.

Olhando para a qualidade de dados fornecidos pelo desafio, tanto para os dados de entrada como para os dados de referência, duas alterações foram realizadas de modo a que os dados se adaptassem à realidade da dissertação.

A primeira alteração foi descartar as características propostas a extrair na segunda fase do desafio, devido ao facto de que a quantidade de características ser insuficiente para caracterizar a informação clínica; a segunda alteração foi realizada sobre o ficheiro de referência da terceira fase, nomeadamente removendo a coluna que indica a existência de ceratose seborreica numa imagem, já que o objetivo desta dissertação é identificar apenas melanoma numa imagem dermatoscópica e não uma das três lesões propostas pelo desafio.

Olhando para os algoritmos de segmentação, com base na seção em que se detalhou esta temática, **todos** os algoritmos mencionados, exceto o **método de *thresholding*** que é um algoritmo de suporte para alguns algoritmos, serão estudados e aplicados para obter a máscara respetiva da lesão da pele segmentada.

Olhando para as características a extrair, com base na seção em que se detalhou esta temática, serão extraídas **todas** as características presentes na Tabela 3, exceto a **área**, o **perímetro**, o **diâmetro maior**, o **diâmetro menor**, a **média** e **desvio-padrão** dos rácios usados na definição da **assimetria da lesão**.

Isto porque as características que fazem parte das propriedades geométricas estão relacionadas com a escala e com as dimensões da imagem dermatoscópicas obtida, não sendo independentes destes dois fatores. As características obtidas para a definição da assimetria da lesão levam um tempo de processamento bastante alto, até 30 minutos ou mais de processamento. Igualmente foi descartado os dados clínicos do paciente, como o sexo e a idade, pois são fatores que foram comprovados ser independentes da existência de melanoma.

Olhando para os algoritmos de aprendizagem automática, com base na seção em que se detalhou esta temática, **todos** os algoritmos propostos serão estudados e aplicados para construir o respetivo modelo de características.

Posto isto, dado que o problema da tese é um problema de aprendizagem automática supervisionada, nomeadamente de classificação, é necessário definir a **fase de treino** (cujo esquema se encontra na Figura 17), a **fase de validação** ou fase de teste (cujo esquema se encontra na Figura 18) e a **fase de aplicação** (cujo esquema se encontra na Figura 19).

Dado que existem 2000 imagens dermatoscópicas, define-se uma boa parte de percentagem de imagens a serem utilizadas na fase de treino (conhecido como dados de treino) e a restante percentagem de imagens a serem utilizadas na fase de validação (conhecido como dados de validação). Tipicamente, 75% ou 80% de dados são utilizados na fase de treino.

Na fase de treino, é aplicado um algoritmo de segmentação para obter as máscaras das lesões das imagens dermatoscópica do conjunto de dados de treino. Uma vez extraídas as características a partir das imagens e das respetivas máscaras, aplica-se um algoritmo de aprendizagem automática para construir um modelo de características.

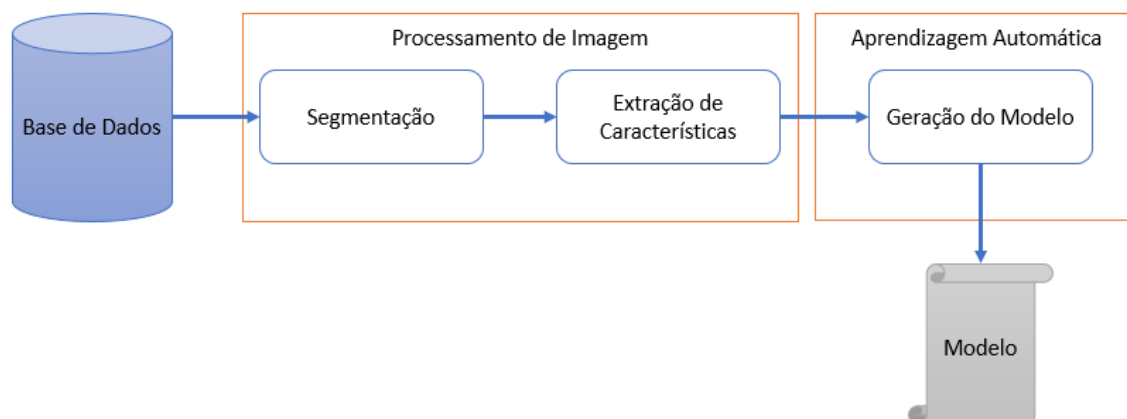


Figura 17: Esquema da Fase de Treino.

Na fase de validação, é aplicado o mesmo algoritmo de segmentação que na fase de treino, para obter as máscaras das lesões das imagens dermatoscópica do conjunto de dados de validação. Uma vez extraídas as características a partir das imagens e das respetivas máscaras, aplica-se o modelo gerado na fase de treino para prever as classificações do conjunto de dados de validação.

As métricas obtidas analisando as classificações reais e as classificações previstas permitiram reanalisar o código, de modo a aperfeiçoar alguma componente das várias implementadas.

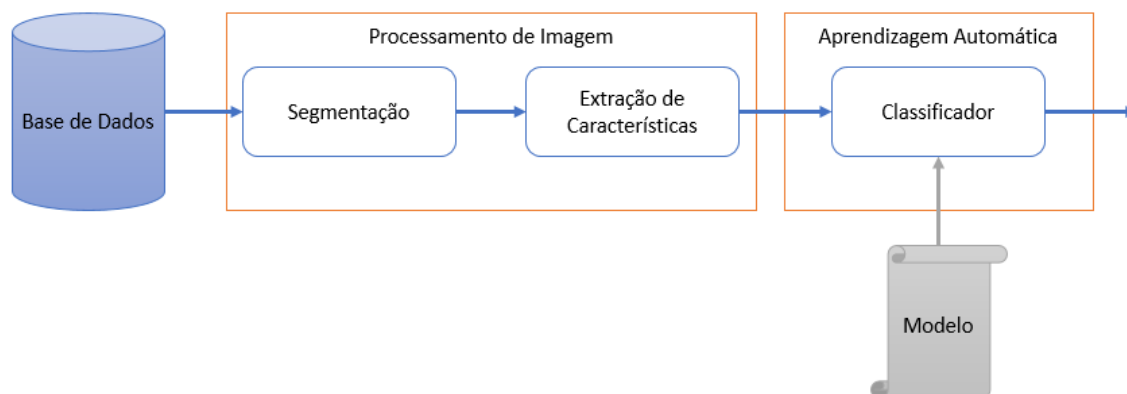


Figura 18: Esquema da Fase de Validação.

Por fim na fase de aplicação, o mesmo algoritmo de segmentação e o mesmo modelo de características das fases anteriores serão utilizadas para identificar melanoma em novas imagens dermatoscópicas.

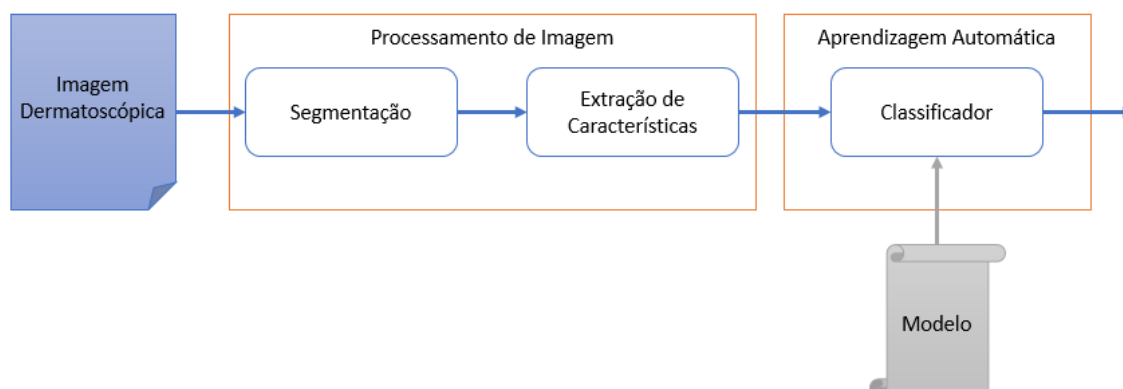


Figura 19: Esquema da Fase de Aplicação.

3.4 Proposta de Validação

Como proposta de validação da solução acima referido, calcular-se-ão métricas para os processos de segmentação e de classificação.

Na segmentação, o objetivo é comparar a máscara obtida por um algoritmo de segmentação (exemplo a máscara representada na Figura 20) com a máscara de referência (exemplo a máscara representada na Figura 21), em que ambas têm a mesma dimensão; isto é, comparar um pixel de uma dada posição numa máscara com o pixel da outra máscara na mesma posição.

Para tal, é contruído uma matriz de confusão para uma segmentação de uma imagem dermatoscópica, onde:

- um **verdadeiro positivo** é um pixel que tem intensidade máxima nas duas máscaras;

- um **verdadeiro negativo** é um pixel que tem intensidade mínima nas duas máscaras;
- um **falso positivo** é um pixel que tem intensidade máxima na máscara obtida, mas na máscara de referência tem intensidade mínima;
- um **falso negativo** é um pixel que tem intensidade mínima na máscara obtida, mas na máscara de referência tem intensidade máxima.

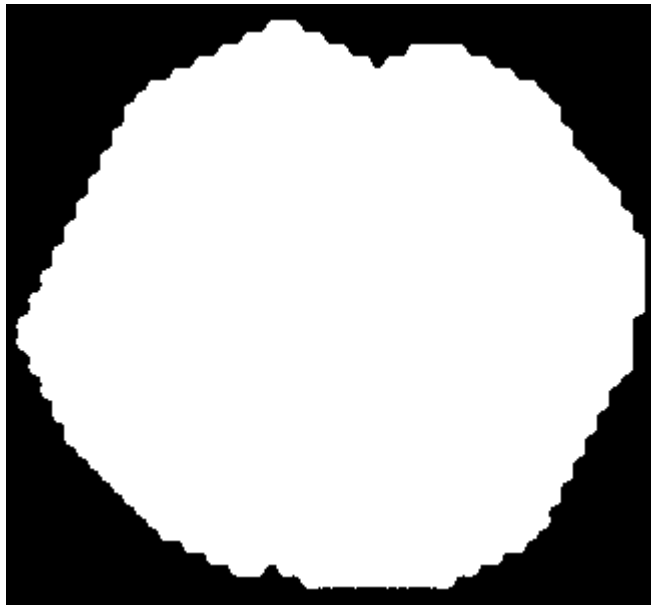


Figura 20: Máscara obtida pelo método de Otsu.

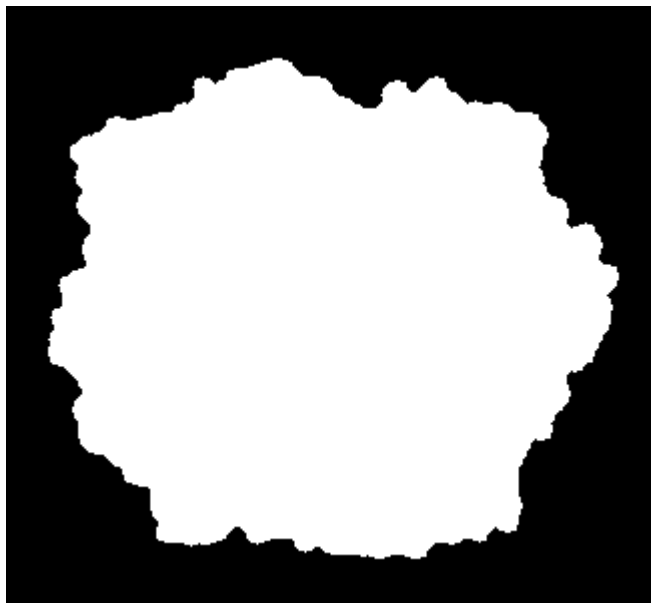


Figura 21: Máscara de referência.

Uma vez contruída a matriz de confusão, aplicar-se-ão cinco métricas, propostas pelo desafio, para validar os vários algoritmos de segmentação. As métricas calculadas, a partir da matriz de confusão, são as seguintes:

- Sensibilidade
- Especificidade
- Exatidão
- Índice de Jaccard
- Índice de Sørensen-Dice

Estas métricas foram calculadas com base numa segmentação de uma imagem dermatoscópica. As métricas para a validação de um algoritmo de segmentação são as médias das métricas de todas as segmentações realizadas em todas as imagens da base de dados.

Na classificação, o objetivo é comparar a classificação que um modelo de características prevê, para um conjunto de características de entrada, com a classificação real deste mesmo conjunto de características. As classificações possíveis, como anteriormente mencionadas, são a existência de melanoma e a não existência de melanoma.

Para tal, é contruído uma matriz de confusão para uma classificação de uma imagem dermatoscópica, onde:

- um **verdadeiro positivo** é uma classificação que indica a existência de melanoma, por parte do classificador, e realmente existe melanoma;
- um **verdadeiro negativo** é uma classificação que indica a não existência de melanoma, por parte do classificador, e realmente não existe melanoma;
- um **falso positivo** é uma classificação que indica a existência de melanoma, por parte do classificador, mas realmente não existe melanoma;
- um **falso negativo** é uma classificação que indica a não existência de melanoma, por parte do classificador, mas realmente existe melanoma.

Uma vez contruída a matriz de confusão, aplicar-se-ão quatro métricas, propostas pelo desafio, para validar os vários algoritmos de segmentação. As métricas calculadas, a partir da matriz de confusão, são as seguintes:

- Sensibilidade
- Especificidade
- Precisão
- Exatidão

•

Capítulo 4

Implementação

4.1 Introdução

Neste capítulo, serão abordados os vários temas relacionados com a implementação do algoritmo proposto, nomeadamente a arquitetura do projeto, a implementação do código relativo à especificação feita no capítulo anterior e a divulgação de resultados.

4.2 Arquitetura do Projeto

4.2.1 Arquitetura Física

Todo o desenvolvimento foi realizado numa máquina **HP Pavilion Sleekbook 15** com o sistema operativo **Windows 10 Home 64-bits** instalado, na versão de compilação 15063.413, num processador **AMD A8-4555M APU** – processador com **quatro núcleos de 64-bit a 1.6 Ghz** de frequência – e **6 GB** de memória RAM.

4.2.2 Arquitetura Tecnológica

O projeto foi desenvolvido principalmente em **Python**, na versão 2.7.10, com o auxílio a bibliotecas/módulos (em inglês, *modules*), que estão identificados na Tabela 5, dado que a linguagem quer os módulos são fáceis de instalar, utilizar e com documentação decente.

Módulo	Versão do Módulo
numpy	1.12.1
opencv	3.2.0
sklearn	0.18.1
skimage	0.11.3

Tabela 5: Versões dos módulos de Python instalados e utilizados no projeto.

Para a instalação da linguagem e respetivos, foi transferido e instalado o software **Python** (**x,y**) já que para além da linguagem em si ser instalada, traz igualmente os módulos apresentados na Tabela 5 instalados, apesar de estar desatualizados. Para a atualização dos módulos, foi utilizado o gestor de módulos do Python denominado **pip**.

Salientar um pormenor importante, de que as matrizes criadas e calculadas pelo módulo **numpy** vem no formato (linha, coluna), enquanto que os restantes módulos tais matrizes vêm no formato (x, y); ou seja, as coordenadas são transpostas entre os módulos mencionados: x corresponde à coluna e y corresponde à linha. Outro pormenor importante é o fato do modelo de cor das imagens tratadas pelo **opencv** for o RGB com a cor vermelha e azul trocadas – isto é, o modelo BGR.

Apesar do software instalado trazer uma interface gráfica de edição e execução, o **spyder**, optou-se por desenvolver na interface gráfica de edição denominada **Notepad++**, na versão 6.9.1, e compilar e executar o código através da linha de comandos do sistema: a **linha de comandos padrão** e, após a grande Atualização dos Criativos de Março de 2017, o **Powershell**.

4.2.3 Estrutura do Projeto

O projeto final relativo à dissertação ficou estruturado da seguinte forma:

- **Source** (em português, código-fonte): pasta contém os ficheiros de código editável e executável.
- **Resources** (em português, recursos): pasta que contém os recursos necessários para a execução do código desenvolvido.
 - **Training Data**: pasta que contém as imagens disponibilizadas pelo desafio.
 - **Segmentation Ground Truth**: pasta que contém as máscaras para cada imagem e que são utilizadas como referência.
 - **Patients Training Data.csv**: ficheiro que contém os dados dos pacientes em que foram obtidas as imagens dermatoscópicas.
 - **Classification Ground Truth.csv**: ficheiro que indica as imagens em que realmente se identificou melanoma ou não.
- **Logs**: pasta que contém resultados obtidos através da execução de código.
 - **Segmentation**: pasta que contém os resultados da segmentação para vários algoritmos; cada resultado de um algoritmo de segmentação está contido numa pasta com o nome do algoritmo em causa.
 - **Segmentation (Validation)**: pasta que contém as métricas calculadas na segmentação a partir das máscaras geradas na pasta **Logs/Segmentation** e as máscaras de referência na pasta **Resources/Segmentation Ground Truth**.
 - **Features**: pasta que contém as características extraídas, em ficheiros CSV.
 - **Classification**: pasta que contém o modelo de características gerado a partir das características da pasta **Logs/Features**.
 - **Classification (Results)**: pasta que contém as classificações previstas pelos modelos da pasta **Logs/Classification**.
 - **Classification (Validation)**: pasta que contém as métricas calculadas na classificação a partir dos modelos na pasta **Logs/Classification** e das classificações da pasta **Logs/Classification (Results)**.

4.3 Desenvolvimento

4.3.1 Solução

Tal como referido no Capítulo 3, a dissertação é composta por duas fases: a fase de modelação (que é composta igualmente por duas fases, a fase de treino e a fase de validação) e a fase de aplicação. Dado a existência de apenas uma base de dados, composta por 2000 imagens e cada imagem associada a si os dados do paciente, uma máscara elaborada por um especialista e pela classificação final, optou-se por desenvolver um conjunto de programas para as três fases em questão.

Na fase de modelação (quer treino quer validação), desenvolveu-se vários programas de modo a facilitar o desenvolvimento. Primeiro, desenvolveu-se um programa **segmentation.py** (cujo esquema se encontra na Figura 22) cujo objetivo era obter a máscara que representa a lesão segmentada de uma imagem dermatoscópica, esta que se encontra na pasta **Resources/Training Data**, através de um conjunto de algoritmos especificados na seção anterior. O resultado da segmentação é guardado na pasta **Logs/Segmentation**.

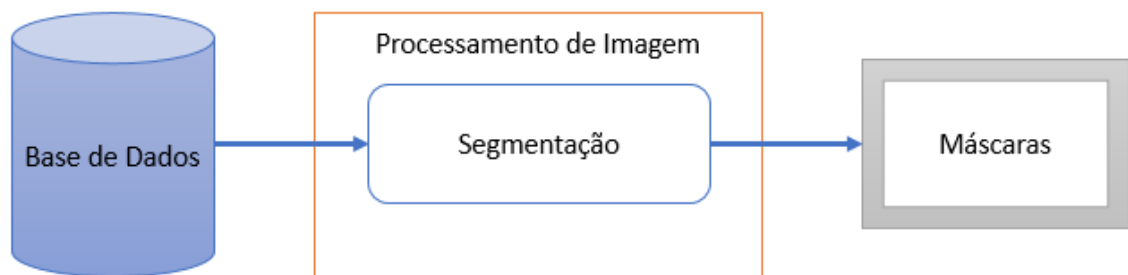


Figura 22: Esquema do programa de segmentação.

De seguida, um programa denominado **extraction.py** (cujo esquema se encontra na Figura 23) permite extrair as características de uma imagem dermatoscópica, especificadas na seção anterior, através da aplicação da respetiva máscara obtida da pasta **Logs/Segmentation**. Estas características são adicionadas a um ficheiro CSV que se encontra na pasta **Logs/Features**.

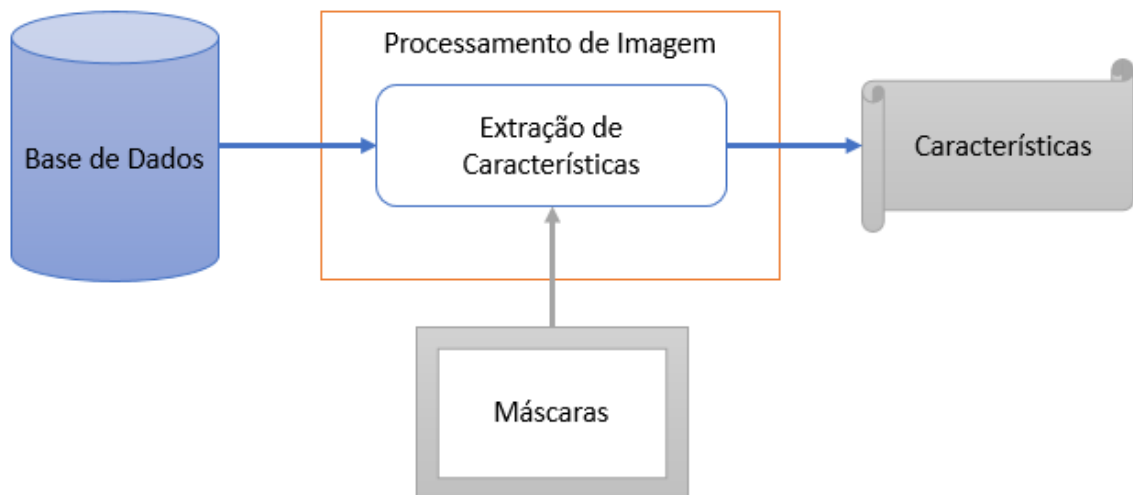


Figura 23: Esquema do programa de extração de características.

Uma vez terminado os programas anteriores na respetiva ordem, dá-se início à fase de treino: o programa **building.py** (cujo esquema se encontra na Figura 24) vai construir um modelo com base no algoritmo de aprendizagem que lhe é fornecido, com base num número de características extraídas com sucesso e que se encontram num ficheiro CSV na pasta **Logs/Features** e vai recorrer às respetivas classificações que se encontram no ficheiro **Resources/Classification Ground Truth.csv**. O resultado obtido pelo algoritmo será um ficheiro PKL, que contém o modelo de aprendizagem automática e será guardado na pasta **Logs/Classification**.

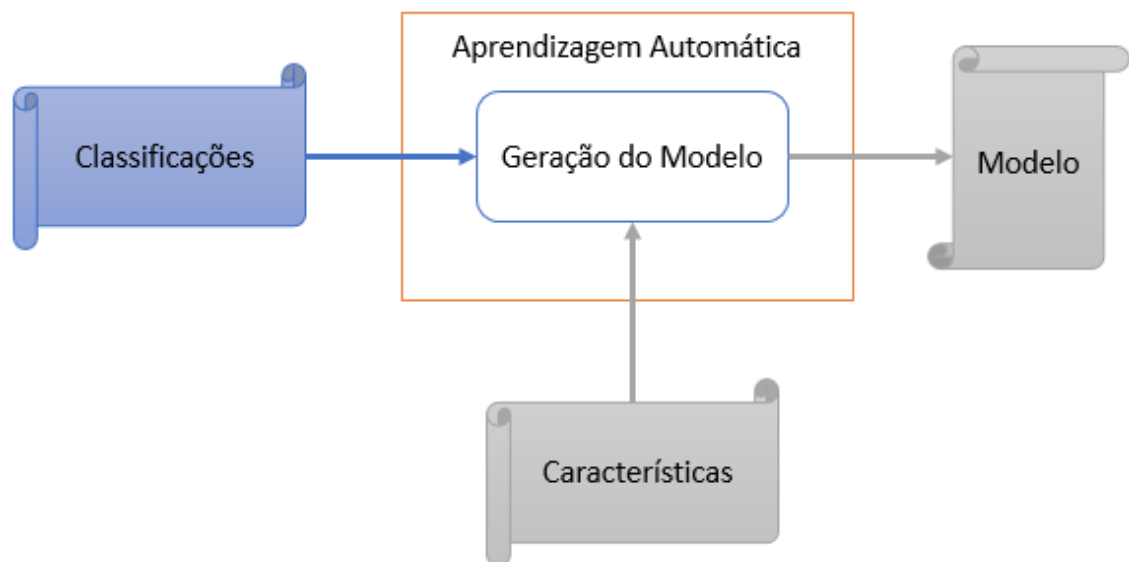


Figura 24: Esquema do programa de construção do modelo.

Na fase de validação, o programa **classification.py** (cujo esquema se encontra na Figura 25) tem como objetivo utilizar os modelos gerados na pasta **Logs/Classification** com o intuito de prever as classificações para um número de características fornecidas da pasta **Logs/Features** e as respetivas classificações reais do ficheiro **Resources/Classification**

Ground Truth.csv. Tais previsões serão guardadas, juntamente com as classificações reais, num ficheiro CSV com o nome do modelo utilizado, na pasta **Logs/Classification (Results)**.

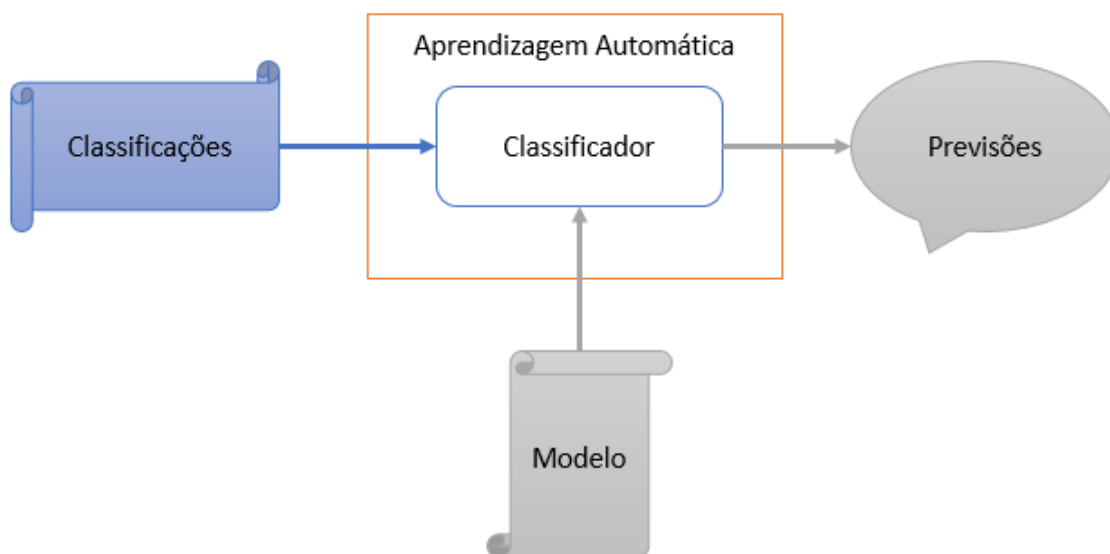


Figura 25: Esquema do programa de previsão de classificações.

Por fim, na fase de aplicação o programa principal **main.py** (cujo esquema se encontra na Figura 19) tem como objetivo utilizar módulos construídos dos programas anteriores para afirmar se existe melanoma ou não numa nova imagem com base num algoritmo de aprendizagem e com base num modelo construído que se encontra na pasta **Logs/Classification**.

4.3.2 Validação

Para validar as soluções implementadas, foram desenvolvidos dois programas, **sMetrics.py** e **cMetrics.py**, de modo a calcular as métricas, especificadas no capítulo anterior, para os resultados dos programas **segmentation.py** e **classification.py**, respetivamente.

O programa **sMetrics.py** irá comparar as máscaras obtidas pelos algoritmos de segmentação e que se encontram na pasta **Logs/Segmentation** com as máscaras de referência que se encontram na pasta **Resources/Segmentation Ground Truth**. As métricas calculadas serão guardadas num ficheiro de texto com o nome do algoritmo de segmentação aplicado na pasta **Logs/Segmentation (Validation)**.

O programa **cMetrics.py** irá comparar as classificações previstas pelos modelos com as classificações reais. Estas duas classificações encontram-se em ficheiros CSV guardadas na pasta **Logs/Classification (Results)**. As métricas calculadas serão guardadas num ficheiro de texto, com o nome do ficheiro CSV utilizado, na pasta **Logs/Classification (Validation)**.

Os resultados obtidos serão detalhados na seguinte seção.

4.4 Resultados

4.4.1 Segmentação

Uma vez calculadas as máscaras para todas as 2000 imagens dermatoscópicas da base de dados, por parte do programa **segmentation.py**, calcularam-se as métricas propostas neste processo, por parte do programa **sMetrics.py**. Tais métricas encontram-se visíveis na Tabela 6.

Método de Segmentação	Tempo Gasto	Métricas (em %)				
		Sensibilidade	Especificação	Exatidão	Índice de Jaccard	Índice de Sørensen-Dice
Otsu	34min	42.48	92.32	85.19	35.73	42.60
K-Means	28h 47min	49.46	89.48	83.57	37.84	44.88
GMM	26h 59min	64.05	60.85	60.04	30.14	37.78
NCut	105h 9min	31.32	94.65	82.59	27.62	31.99
Watershed	10h 39min	45.66	79.97	76.63	29.59	36.74
Closed Shortest Path	1h 6min	81.64	93.88	89.29	65.25	75.03

Tabela 6: Métricas obtidas pelos métodos de segmentação.

Pode-se concluir que de forma geral, o algoritmo *Closed Shortest Path* apresenta melhores resultados, pois em todas as métricas exceto a especificação apresenta o maior valor. Sendo assim, este algoritmo será utilizado como padrão na obtenção da máscara da lesão.

4.4.2 Classificação

Para este processo, construiu-se um modelo de características para cada um dos cinco algoritmos de aprendizagem automática e para uma das duas bases de segmentação; isto é, um modelo de características utilizando as máscaras de referência e um modelo de características utilizando as máscaras obtidas pelo algoritmo *Closed Shortest Path*, na extração de características pelo programa **extraction.py**.

A utilização das máscaras de referência nesta fase tem como objetivo verificar a qualidade das características extraídas, sendo que os modelos de características gerados não vão ser utilizados na identificação automática de melanoma. Foi assim criado dois ficheiros CSV com as características extraídas, um para cada base de segmentação.

No programa **building.py**, foram utilizados 75% das imagens da base de dados como dados de treino; isto é, 1500 imagens dermatoscópicas. As restantes 25% de imagens foram utilizadas como dados de validação; isto é, 500 imagens dermatoscópicas. Salientar igualmente que 135 das 500 imagens dermatoscópicas de validação estão identificadas melanoma.

Uma vez gerado os dez modelos de características acima referido, utilizou-se o programa **classification.py** para prever que classificação cada modelo previa para as respetivas imagens dos dados de validação. Uma vez previstos as classificações, o programa **cMetrics.py** calculou as métricas que se encontram visíveis na Tabela 7 e Tabela 8.

Base de Segmentação	Algoritmo de Aprendizagem Automática	Matriz de Confusão			
		VP	FP	FN	VN
Referência	DT	22	79	113	286
	ANN	0	0	135	365
	SVM	0	0	135	365
	RF	6	18	129	347
	NB	56	129	79	236
<i>Closed Shortest Path</i>	DT	31	94	104	271
	ANN	0	0	135	365
	SVM	0	0	135	365
	RF	13	43	122	322
	NB	0	0	135	365

Tabela 7: Matriz de Confusão para um modelo de cada base de segmentação.

Base de Segmentação	Algoritmo de Aprendizagem Automática	Métricas (em %)			
		Sensibilidade	Especificação	Precisão	Exatidão
Referência	DT	16.30	78.36	21.78	61.6
	ANN	0	100	-	73.00
	SVM	0	100	-	73.00
	RF	4.44	95.07	25.00	70.6
	NB	41.48	64.66	30.27	58.4
<i>Closed Shortest Path</i>	DT	22.96	74.25	24.80	60.40
	ANN	0	100	-	73.00
	SVM	0	100	-	73.00
	RF	9.63	88.22	23.21	67.00
	NB	0	100	-	73.00

Tabela 8: Métricas calculadas para cada modelo de uma base de segmentação.

Analisando as métricas calculadas, pode-se afirmar que a existência de um grande número de falsos negativos é bastante preocupante. Tanto as **máquinas de suporte vetorial** como as **redes neurais artificiais** não permitiram obter qualquer verdadeiro positivo; isto é, não conseguiram prever imagens em que estava realmente identificado melanoma. Igualmente se verifica o mesmo caso para o algoritmo **Naive Bayes** quando se utilizam as máscaras das lesões obtidas pelo algoritmo *Closed Shortest Path*.

Sendo assim, opta-se em utilizar o modelo de características geradas pela **Floresta Aleatória**, pois entre esse modelo e o modelo obtido pela **Árvore de Decisão**, o primeiro modelo apresenta uma maior exatidão, apesar de ter uma menor sensibilidade.

Não se pode descartar que as exatidões na base de segmentação de referência, numa maneira geral, foram superiores às exatidões na base de segmentação do algoritmo *Closed Shortest Path*, igualmente sendo possível, na base de segmentação de referência, prever verdadeiros positivos no caso do algoritmo de Naive Bayes. Apesar de tudo, pode-se igualmente afirmar que as sensibilidades na base de segmentação do algoritmo *Closed Shortest Path* foram ligeiramente superiores às sensibilidades na base de segmentação de referência.

Significa que as máscaras de referência permitem classificar melhor do que as máscaras obtidas pelo algoritmo *Closed Shortest Path*.

De forma geral, as métricas obtidas permitem afirmar que as características extraídas não representam adequadamente as imagens dermatoscópicas (ou têm baixa qualidade), dado que para os modelos que previram pelo menos verdadeiros positivos, apresentam uma exatidão abaixo dos 75% e uma sensibilidade abaixo dos 50%, o que fica bastante a desejar.

Capítulo 5

Conclusão

O objetivo claro de identificar melanoma de forma automática, com base em imagens dermatoscópicas revelou-se bastante complexo para a ideia que tinha inicialmente da dissertação.

Tendo em conta o tempo dado para o desenvolvimento da dissertação era bastante curto, o objetivo foi cumprido – foi possível identificar melanoma de forma automática em certas imagens dermatoscópicas. Os resultados obtidos no processo de segmentação das lesões foram agradáveis, mas os resultados obtidos através dos modelos de características ficaram aquém das expectativas, pois era esperado que pelo menos um modelo gerado tivesse resultados com uma exatidão acima dos 75%, o que significaria que o algoritmo classificaria corretamente $\frac{3}{4}$ das imagens utilizadas na fase de validação.

Várias dificuldades foram sentidas. Em questões técnicas, aconteceram várias adversidades, nomeadamente problemas de versões dos módulos de Python, em que várias funções de módulos ou não existiam, de uma versão anterior para mais recente e vice-versa, ou apresentavam erros ao serem invocadas. Outra adversidade foi a atualização dos criativos de Março de 2017 para o Windows 10, que desinstalou o Python e os vários módulos como também removeu vários drivers do computador em que se desenvolveu esta dissertação.

Em questões profissionais, dada a minha inexperiência nas áreas em que a dissertação se enquadrava, bastante tempo foi gasto em aprender conceitos e técnicas que iria ser aplicadas nas várias componentes do trabalho para alcançar o objetivo.

Como trabalho futuro, seria importante investir mais na qualidade das características, pois são estas as que são utilizadas na construção do modelo. Um investimento em melhorar o desempenho da obtenção de características relativas à assimetria de uma lesão seria importante. Igualmente investir noutros modelos de cor que não seja o RGB e o HSV seria importante de experimentar como também introduzir novas características relativamente à textura que não seja

só a matriz de coocorrências seria desejável. Um investimento sério no pré-processamento da imagem – tal como remoção de pelos ou normalização da imagem – pode ser feito.

Referências

- [1] Euromelanoma, “Melanoma”, 17 de Abril de 2017, <http://www.euromelanoma.org/intl/node/77>
- [2] World Cancer Research Fund International, “Worldwide Data”, 17 de Abril de 2017, <http://www.wcrf.org/int/cancer-facts-figures/worldwide-data>
- [3] International Skin Imaging Collaboration, “ISIC 2017: Skin Lesion Analysis Towards Melanoma Detection”, 17 de Abril de 2017, https://challenge.kitware.com/#challenge/n/ISIC_2017%3A_Skin_Lesion_Analysis_Towards_Melanoma_Detection
- [4] ISBI, “About ISBI”, 17 de Abril de 2017, <http://biomedicalimaging.org/2017/contact-us/>
- [5] International Society for Digital Imaging of the Skin, “ISIC Project”, 17 de Abril de 2017, <http://www.isdis.net/index.php/isic-project>
- [6] WebMD, “Picture of the Skin”, 19 de Abril de 2017, <http://www.webmd.com/skin-problems-and-treatments/picture-of-the-skin#1>
- [7] Wikipédia, 14 de Abril de 2017, “Melanócito”, 19 de Abril de 2017, <https://pt.wikipedia.org/wiki/Melan%C3%B3cito>
- [8] Wikipédia, 12 de Janeiro de 2017, “Neoplasm”, 17 de Janeiro de 2017, <https://en.wikipedia.org/wiki/Neoplasm>
- [9] SkinVision, “What causes skin lesions?”, 19 de Abril de 2017, <https://skinvision.com/library/what-causes-skin-lesions>

- [10] American Academy of Dermatology, “Seborrheic keratoses: Overview”, 19 de Abril de 2017, <http://emedicine.medscape.com/article/1059477-overview>
- [11] Medscape, “Seborrheic Keratosis: Background”, 19 de Abril de 2017, <https://www.aad.org/public/diseases/bumps-and-growths/seborrheic-keratoses>
- [12] PRAXISDIENST, “Dermatoscopes”, 19 de Abril de 2017, <https://www.praxisdienst.com/en/Human/Diagnostics/Specialised+Diagnostics/Dermatoscopes/>
- [13] molemap, “HOW TO IDENTIFY MELANOMA”, 19 de Abril de 2017, <https://molemap.co.nz/melanoma-skin-cancer/how-to-identify-melanoma/>
- [14] Johr, R. H.; Junho de 2002; “Dermoscopy: alternative melanocytic algorithms—the ABCD rule of dermatoscopy, menzies scoring method, and 7-point checklist”; 19 de Abril de 2017
- [15] Gio Winderhold, John McCarthy; Maio de 1992; “Arthur Samuel: pioneer in machine learning”; IBM Journal of Research and Development; 36º Volume; 3º Exemplar, Páginas 328-331
- [16] Thomas M. Mitchell; 1997; “Machine Learning”; McGraw-Hill, Inc.; Nova Iorque, NY, Estados Unidos da América
- [17] J. Ross Quinlan; 1993; "Induction of decision trees"; Morgan Kaufmann Publishers, Inc.; São Francisco, CA, Estados Unidos da América
- [18] J. Ross Quinlan; 1993; "C4.5 Programs for Machine Learning"; Morgan Kaufmann Publishers, Inc.; São Francisco, CA, Estados Unidos da América
- [19] College of Engineering of Univerity of Florida, “The ID3 Algorithm”, 24 de Abril de 2017, <https://www.cise.ufl.edu/~ddd/cap6635/Fall-97/Short-papers/2.htm>
- [20] Dr. Saed Sayad, “Decision Tree - Classification”, 24 de Abril de 2017, http://www.saedsayad.com/decision_tree.htm
- [21] Leo Breiman, Jerome Friedman, Charles J. Stone, R.A. Olshen; 1984; “Classification and regression trees”; Wadsworth & Brooks/Cole Advanced Books & Software; Monterey, CA, Estados Unidos da América
- [22] Wikipédia, 17 de Abril de 2017, “Artificial neural network”, 25 de Abril de 2017, https://en.wikipedia.org/wiki/Artificial_neural_network
- [23] Warren McCulloch; Walter Pitts; 1943; "A Logical Calculus of Ideas Immanent in Nervous Activity"; Bulletin of Mathematical Biophysics
- [24] Frank Rosenblatt; 1957; “The Perceptron--a perceiving and recognizing automaton”; Cornell Aeronautical Laboratory

- [25] Corinna Cortes, Vladimir Vapnik; Setembro de 1995; "Support-vector networks"; AT&T Bell Labs.; 20º Volume; 3ª Edição, Páginas 273–297
- [26] Bernhard E. Boser, Isabelle M. Guyon, Vladimir N. Vapnik; Julho de 1992; “A training algorithm for optimal margin classifiers”; COLT '92 Proceedings of the fifth annual workshop on Computational learning theory; Nova Iorque, NY, Estados Unidos da América
- [27] KDnuggets, Julho de 2016, “Support Vector Machines: A Simple Explanation”, 25 de Abril de 2017, <http://www.kdnuggets.com/2016/07/support-vector-machines-simple-explanation.html>
- [28] Quora, 2014, “What does support vector machine (SVM) mean in layman's terms?”, 25 de Abril de 2017, <https://www.quora.com/What-does-support-vector-machine-SVM-mean-in-laymans-terms>
- [29] Wikipédia, 17 de Abril de 2017, “Support vector machine”, 25 de Abril de 2017, https://en.wikipedia.org/wiki/Support_vector_machine#/media/File:Kernel_Machine.png
- [30] Tin Kam Ho; Agosto de 1995; “Random Decision Forests”; Proceedings of the 3rd International Conference on Document Analysis and Recognition; Montreal, QC, Estados Unidos da América
- [31] YouTube, 4 de Abril de 2014, “How Random Forest algorithm works”, 26 de Abril de 2017, <https://www.youtube.com/watch?v=loNcrMjYh64>
- [32] YouTube, 17 de Junho de 2016, “Random Forest based Classification”, 26 de Abril de 2017, <https://www.youtube.com/watch?v=ajTc5y3OqSQ>
- [33] Stuart Russell, Peter Norvig; 1995; “Artificial Intelligence: A Modern Approach”; Prentice Hall; 1º Volume; 2ª Edição
- [34] Harold Jeffreys; 1973; “Scientific Inference”; Cambridge University Press; 1º Volume; 3ª Edição
- [35] Analytics Vidhya, 13 de Setembro de 2015, “6 Easy Steps to Learn Naive Bayes Algorithm (with code in Python)”, 26 de Abril de 2017, <https://www.analyticsvidhya.com/blog/2015/09/naive-bayes-explained/>
- [36] YouTube, 27 de Agosto de 2014, “How Naive Bayes Classifier Works 1/2.. Understanding Naive Bayes and Example”, 26 de Abril de 2017, <https://www.youtube.com/watch?v=XcwH9JGfZOU>
- [37] J. B. MacQueen; 1967; “Some Methods for classification and Analysis of Multivariate Observations”; Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability; University of California Press

- [38] YouTube, 19 de Janeiro de 2014, “K-means clustering: how it works”, 26 de Abril de 2017, https://www.youtube.com/watch?v=_aWzGGNrcic
- [39] B.S. Everitt, D.J. Hand; 1981; “Finite mixture distributions”; Chapman & Hall
- [40] Quora, 2014, “What is the difference between K-means and the mixture model of Gaussian?”, 27 de Abril de 2017, <https://www.quora.com/What-is-the-difference-between-K-means-and-the-mixture-model-of-Gaussian>
- [41] Rahman Farnoosh, Gholamhossein Yari, Behnam Zarpak; “Image Segmentation using Gaussian Mixture Models”,
http://djafari.free.fr/maxent2006/Papers/008_Zarpak/008_Zarpak.pdf
- [42] Mohand Saïd Allili; 2010; “A short tutorial on Gaussian Mixture Models”; Université du Québec, Outaouais;
http://www.computerrobotvision.org/2010/tutorial_day/GMM_said_crv10_tutorial.pdf
- [43] BRILLIANT, “Gaussian Mixture Model”, 27 de Abril de 2017, <https://brilliant.org/wiki/gaussian-mixture-model/>
- [44] A.P. Dempster, N.M. Laird; 1977; "Maximum Likelihood from Incomplete Data via the EM Algorithm"; Journal of the Royal Statistical Society; Série B; 39º Volume; Páginas 1–38
- [45] Linda G. Shapiro, George C. Stockman; 2001; “Computer Vision”; Prentice-Hall; New Jersey, NY, Estados Unidos da América
- [46] Robert Hirsch; 2004; “Exploring Colour Photography: A Complete Guide”; Laurence King Publishing
- [47] George H. Joblove, Donald Greenberg; Agosto de 1978; "Color spaces for computer graphics". Série Computer Graphics; 12º Volume; Edição 3: Páginas 20–25
- [48] OpenCV Documentation, “Miscellaneous Image Transformations”, 2 de Maio de 2017,
http://docs.opencv.org/2.4/modules/imgproc/doc/miscellaneous_transformations.html?highlight=threshold
- [49] Quora, 2015, “What are the roles of the threshold function and of a sigmoid function on artificial neural networks?”, 2 de Maio de 2017, <https://www.quora.com/What-are-the-roles-of-the-threshold-function-and-of-a-sigmoid-function-on-artificial-neural-networks>
- [50] Nobuyuki Otsu; 1979; "A threshold selection method from gray-level histograms"; IEEE Trans. Sys., Man., Cyber; 9º Volume; Páginas 62–66

- [51] OpenCV Documentation, “Image Thresholding”, 2 de Maio de 2017, http://docs.opencv.org/trunk/d7/d4d/tutorial_py_thresholding.html
- [52] Fernand Meyer; 1991; “Un algorithme optimal pour la ligne de partage des eaux”; 8me congrès de reconnaissance des formes et intelligence artificielle; 2º Volume; Páginas 847–857; Léon, França
- [53] L. Najman, M. Schmitt; 1994; “Watershed of a continuous function. In Signal Processing (Special issue on Mathematical Morphology.)”; 38º Volume; Páginas 99–112
- [54] M. Kass, A. Witkin, D. Terzopoulos; 1988; "Snakes: Active contour models"; International Journal of Computer Vision; 1º Volume; 4º Edição
- [55] David Karger; 1993; "Global Min-cuts in RNC and Other Ramifications of a Simple Mincut Algorithm"; Proc. 4th Annual ACM-SIAM Symposium on Discrete Algorithms
- [56] Stoer Wagner, Mechthild Wagner, Frank Wagner; 1997; "A simple min-cut algorithm"; Journal of the ACM; 44º Volume; 4º Edição; Páginas 585–591.
- [57] Jianbo Shi, Jitendra Malik; 1997; "Normalized Cuts and Image Segmentation"; IEEE Conference on Computer Vision and Pattern Recognition; Páginas 731–737
- [58] YouTube, 23 de Fevereiro de 2015, “Measure Affinity”, 10 de Maio de 2017, <https://www.youtube.com/watch?v=zyOtB7rffyE>
- [59] CODEFORCES, 2015, “DFS explanation for beginners”, 15 de Maio de 2017, <http://codeforces.com/blog/entry/16823>
- [60] Arthur B. Kahn; 1962; "Topological sorting of large networks", Communications of the ACM, 5º Volume; 11º Edição; Páginas 558–562
- [61] Jaime S. Cardoso, Inês Domingues, Hélder P. Oliveira; Agosto de 2014; “Closed Shortest Path in the Original Coordinates with an Application to Breast Cancer”, 28 de Abril de 2017, <http://www.inescporto.pt/~jsc/publications/journals/2014JaimeIJPRAI.pdf>
- [62] OpenCV Documentation, “Contours in OpenCV”, 18 de Maio de 2017, http://docs.opencv.org/trunk/d3/d05/tutorial_py_table_of_contents_contours.html
- [63] Robert M Haralick, K Shanmugam, Its'hak Dinstein; 1973; "Textural Features for Image Classification"; IEEE Transactions on Systems, Man, and Cybernetics. 3º Volume; 6º Edição; Páginas 610–621
- [64] scikit-image Documentation, “skimage.feature”, 18 de Maio de 2017, <http://scikit-image.org/docs/dev/api/skimage.feature.html#skimage.feature.greycomprops>

- [65] M. K. Hu; 1962; "Visual Pattern Recognition by Moment Invariants", IRE Trans. Info. Theory; 8º Volume; Páginas 179–187
- [66] Tom Fawcett; 2006; "An Introduction to ROC Analysis"; Pattern Recognition Letters; 27º Volume; 8º Edição; Páginas 861–874
- [67] Paul Jaccard; 1901; "Étude comparative de la distribution florale dans une portion des Alpes et des Jura"; Bulletin de la Société Vaudoise des Sciences Naturelles; 37º Volume; Páginas 547–579.
- [68] Paul Jaccard; 1912; "The distribution of the flora in the alpine zone"; New Phytologist; 11º Volume; Páginas 37–50
- [69] Lee R. Rice; 1945; "Measures of the Amount of Ecologic Association Between Species"; Ecology; 26º Volume; 3º Edição; Páginas 297–302
- [70] T. Sørensen; 1948; "A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons"; Kongelige Danske Videnskabernes Selskab; 5º Volume; 4º Edição; Páginas 1–34.
- [71] Margarida Silveira, Jacinto C. Nascimento, Jorge S. Marques, André R. S. Marçal, Teresa Mendonça, Syogo Yamauchi; Fevereiro de 2009; "Comparison of Segmentation Methods for Melanoma Diagnosis in Dermoscopy Images"; IEEE JOURNAL OF SELECTED TOPICS IN SIGNAL PROCESSING; 3º Volume
- [72] Roberta Barbosa Oliveira; Março de 2017; "Pattern Recognition in Pigmented Skin Lesion Images using Ensemble Methods"; Dissertação; Faculdade de Engenharia da Universidade do Porto
- [73] Ana Catarina Fidalgo Barata, Fevereiro de 2017; "Automatic Detection of Melanomas Using Dermoscopy Images"; Doutoramento; Instituto Superior Técnico da Universidade de Lisboa
- [74] Munya A. Arasi, El-Sayed A. El-Dahshan, El-Sayed M. El-Horbaty, AbdelBadeeh M. Salem; 3 de Setembro de 2016; "Malignant Melanoma Detection Based on Machine Learning Techniques: A Survey"; Dept. of Computer Science, Faculty of Computer and Information Sciences, Ain Shams University, Abbassia, Cairo, Egipto